



Európska únia
Európsky fond regionálneho rozvoja



Detekcia komunit v komplexných sieťach

Marek Ciglan a Michal Laclavík,
Ústav informatiky SAV,

marek.ciglan@savba.sk, michal.laclavik@savba.sk



Kompetenčné centrum inteligentných technológií pre elektronizáciu a informatizáciu systémov a služieb

KC-INTELINSYS, ITMS: 26240220072

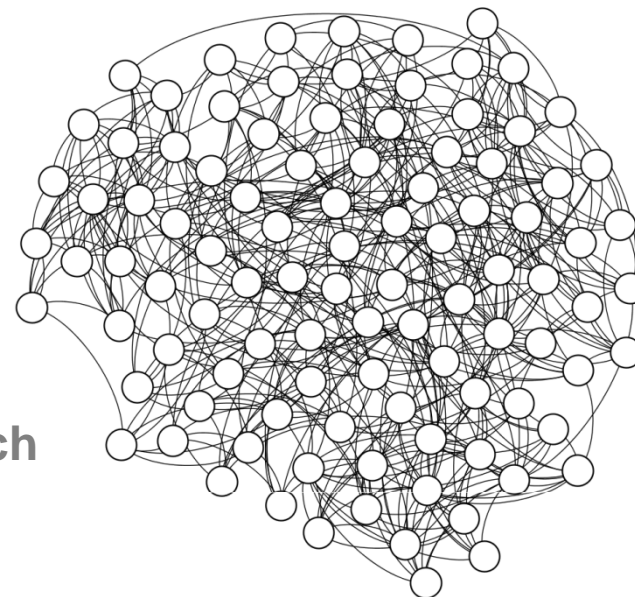


Podporujeme výskumné aktivity na Slovensku/Projekt je spolufinancovaný zo zdrojov EÚ.

Na realizáciu projektu sme získali nenávratný finančný príspevok v rámci Operačného programu výskum a vývoj.



- Úvod
 - Komplexné siete v reálnom svete
 - Vlastnosti komplexných sietí
 - Ako vyzerajú?
 - Porovnanie s náhodným sieťami.
 - Typické úlohy dolovania v grafových dátach
- Detekcia komunít
 - Definícia úlohy
 - Neexistencia konsenzu na formalizácii úlohy a dôsledky
 - Najrozšírenejšie špecifikácie problému používané v literatúre
 - Triedy algoritmov na detekciu komunít
- Meranie kvality zhlukovania
 - Lokálne metriky kvality komunít
 - Modularita





- Meranie kvality zhlukovania
 - **Moudlarita**
 - Popis
 - Problémy spojené s modularitou
 - **Benchmarky pre detekciu komunit**
 - GN benchmark
 - LFR benchmark
 - **Porovnanie kanonického a detekovaného rozdelenia**
 - NMI
 - Priemerná Jaccardová podobnosť
- Modifikácie základnej úlohy
 - **Perekrývajúce sa zhluky**
 - **Vrcholy / hrany s atribútmi**
 - **Detekcia komunity pre daný vrchol**



- Rýchle algoritmy na detekciu komunit
 - Label propagation
 - Louvain method
 - SCCD
- Zhlukovanie na sieťach reálneho sveta
 - Ako dobre fungujú algoritmy na detekciu komunit na sieťach reálneho sveta s explicitne definovanými komunitami

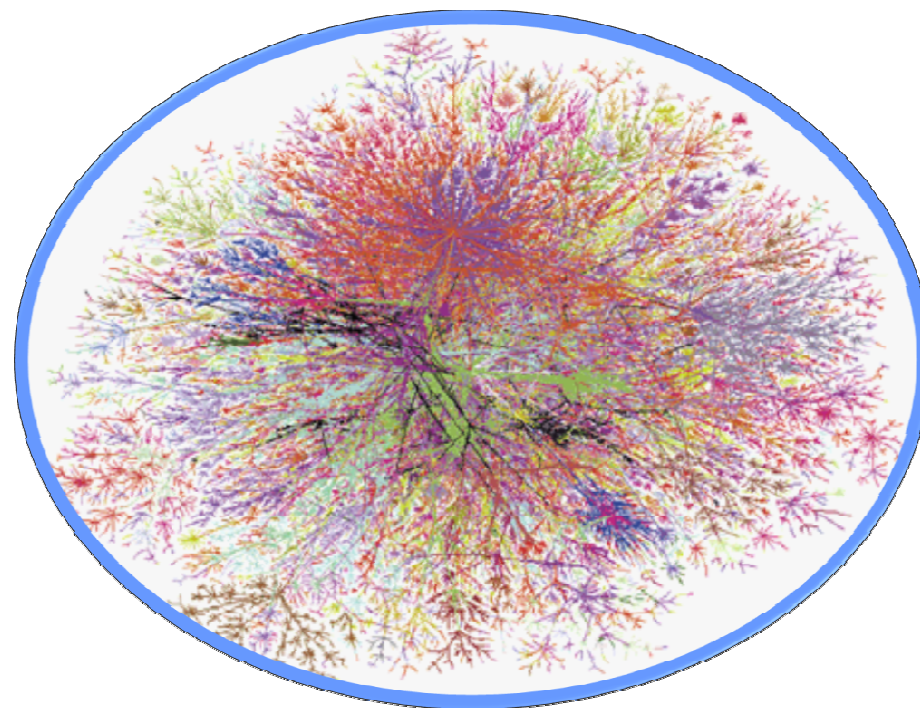


Európska únia
Európsky fond regionálneho rozvoja

Komplexné siete



- **Grafové dáta** – zachytávajú entity (vrcholy grafu) a ich vzťahy (hrany grafu)
- **Veľké dáta** – rozsah spracovávaných dát robí orientáciu v dátach, ich vizualizáciu a štúdium ich vlastností netriviálnym problémom
- **Komplexný systém** – zložený zo spojených častí, ako celok vykazuje vlastnosti, ktoré sa nedajú pozorovať na úrovni jednotlivých častí





- Čoraz viac dostupných dát, ktoré možno prirodzene modelovať ako komplexné siete
 - **Veľké on-line aplikácie:**
 - Sociálne siete (facebook, twitter)
 - Komunikácia (instant messaging, emaily, call networks)
 - Sociálne média (blogy)





- Sociálna sieť – relácie medzi užívateľmi
 - Facebook - 750 miliónov užívateľov
 - YouTube - 490 miliónov pravidelných užívateľov
 - Twitter - 550 miliónov užívateľov
 - Wikipedia - 91,000 kontribútorov
- Obsah generovaný užívateľmi
 - textové správy, fotky, videá, reakcie (+1 / Likes)
 - Facebook - 30 miliárd zdieľaných položiek / mesiac
 - Twitter - 190 miliónov mikropostov / deň
 - Wikipedia - 17 miliónov hostovaných článkov
- Interakcia užívateľov s obsahom
 - YouTube - 92 miliárd zobrazení stránok / mesiac
 - Twitter - 1.6 miliárd dopytov za deň
 - Facebook – priemerný čas strávený na stránke za mesiac: 15H 33M



- Sociálne siete
 - On-line sociálne siete
 - Komunikačné siete
- Informačné siete
 - Blogy
 - Citačné siete
 - WWW, hypertext
- Sémantické siete
 - Linked open data cloud
- Jazykové siete
 - Term co-occurrence networks
- Technologické siete
 - Cestné siete
 - Inžinierske siete (elektrické, potrubné)



Európska únia
Európsky fond regionálneho rozvoja

Vlastnosti komplexných sietí



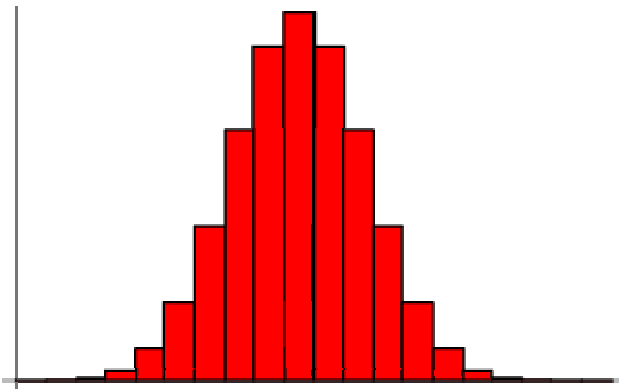
- Ako vyzerajú komplexné siete z reálneho sveta?
- Ako sa podobajú reálne siete na náhodné grafy?
- Majú rovnaké typy sietí podobné vlastnosti?
 - **Sú globálne vlastnosti sociálnej siete A podobné vlastnostiam sociálnej siete B?**



- Podobajú sa komplexné siete na náhodné grafy?
- Erdos-Renyi model náhodného grafu:
 - **Variant 1:**
 - Graf o n vrcholoch, každá hrana (i,j) existuje s pravdepodobnosťou p
 - Teda graf s m hranami sa vyskytuje s pravdepodobnosťou:
 - $p^m * (1-p)^{M-m}$; kde $M=n(n-1)/2$
 - **Variant 2:**
 - Graf on n vrcholoch a m náhodne vybraných hranách



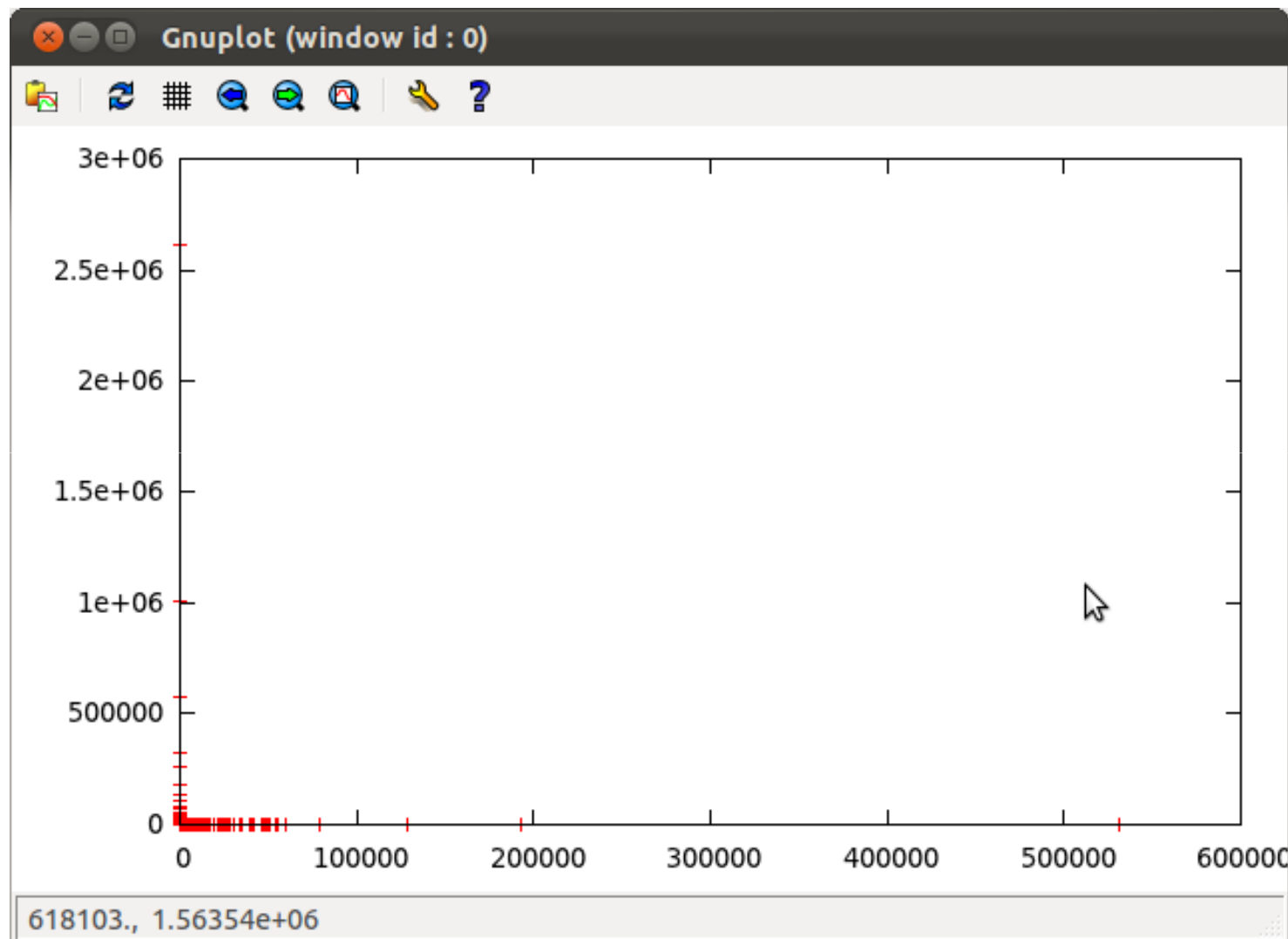
- Podobajú sa komplexné siete na náhodné grafy?
- Erdos-Renyi model náhodného grafu:
 - **Variant 1:**
 - Graf o n vrcholoch, každá hrana (i,j) existuje s pravdepodobnosťou p
 - Teda graf s m hranami sa vyskytuje s pravdepodobnosťou:
 - $p^m * (1-p)^{(M-m)}$; kde $M=n(n-1)/2$
 - **Variant 2:**
 - Graf on n vrcholoch a m náhodne vybraných hranách
- Distribúcia stupňov náhodného grafu - binomická



$$p_k = \binom{n}{k} p^k (1-p)^{n-k}$$



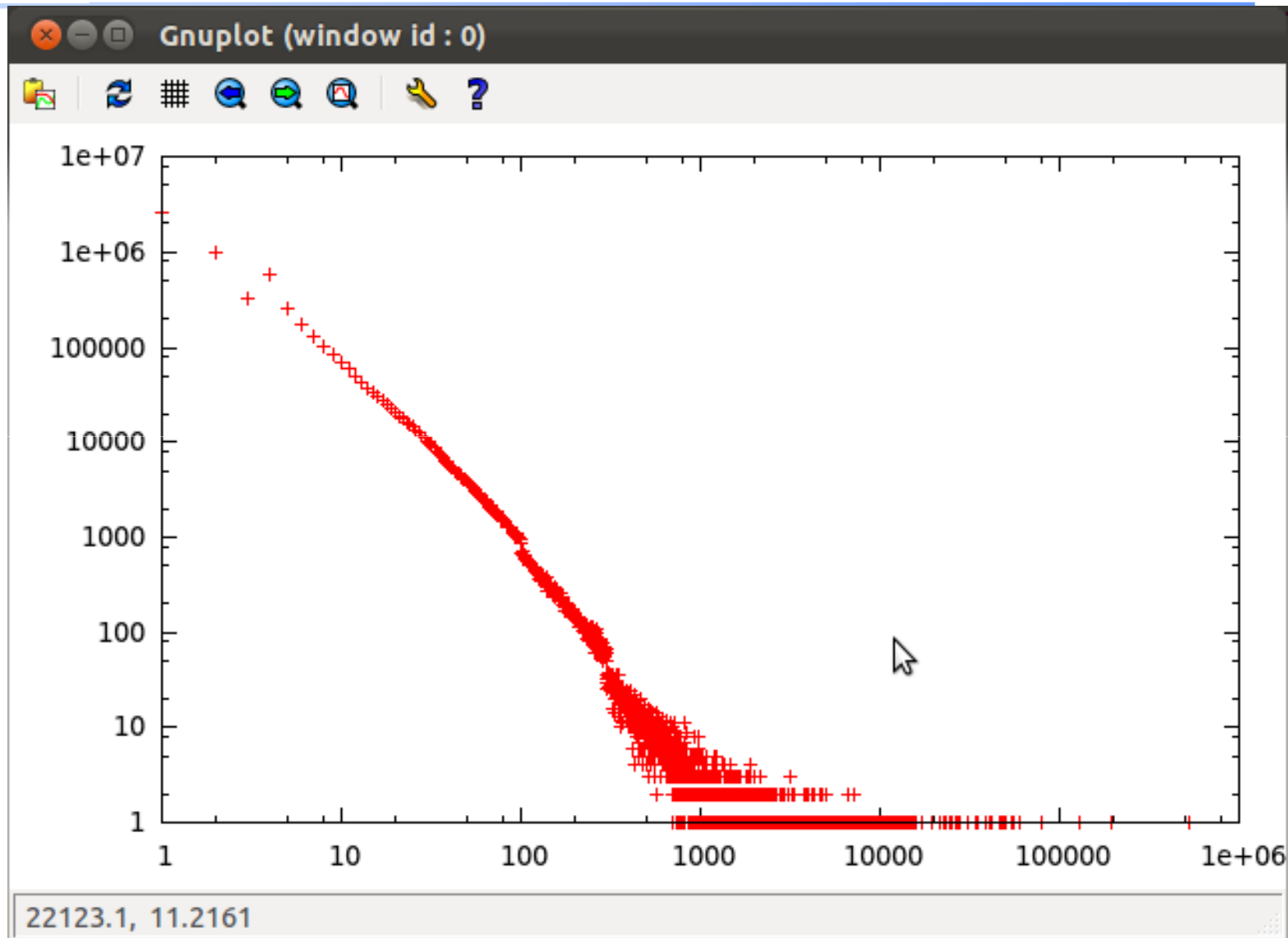
- Podobajú sa komplexné siete na náhodné grafy?
- Erdos-Renyi model náhodného grafu:
 - **Variant 1:**
 - Graf o n vrcholoch, každá hrana (i,j) existuje s pravdepodobnosťou p
 - Teda graf s m hranami sa vyskytuje s pravdepodobnosťou:
 - $p^m * (1-p)^{M-m}$; kde $M=n(n-1)/2$
 - **Variant 2:**
 - Graf o n vrcholoch a m náhodne vybraných hranách
- Distribúcia stupňov náhodného grafu – binomická
- Priemer grafu: $O(\log n)$ - zväčšujúci sa s veľkosťou grafu





Európska únia
Európsky fond regionálneho rozvoja

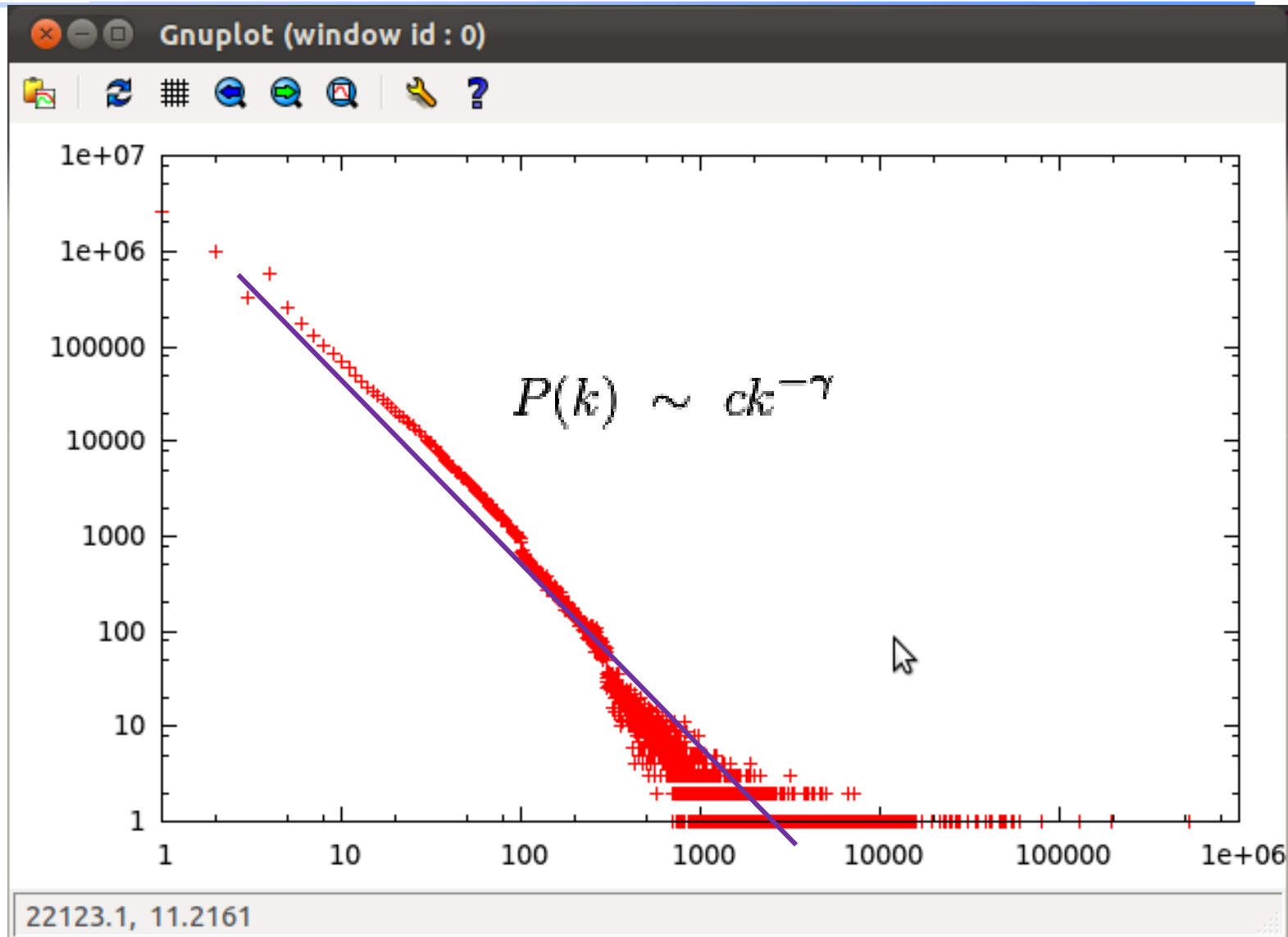
Distribúcia stupňov v reálnych sieťach - log škála (príklad DBpedia)

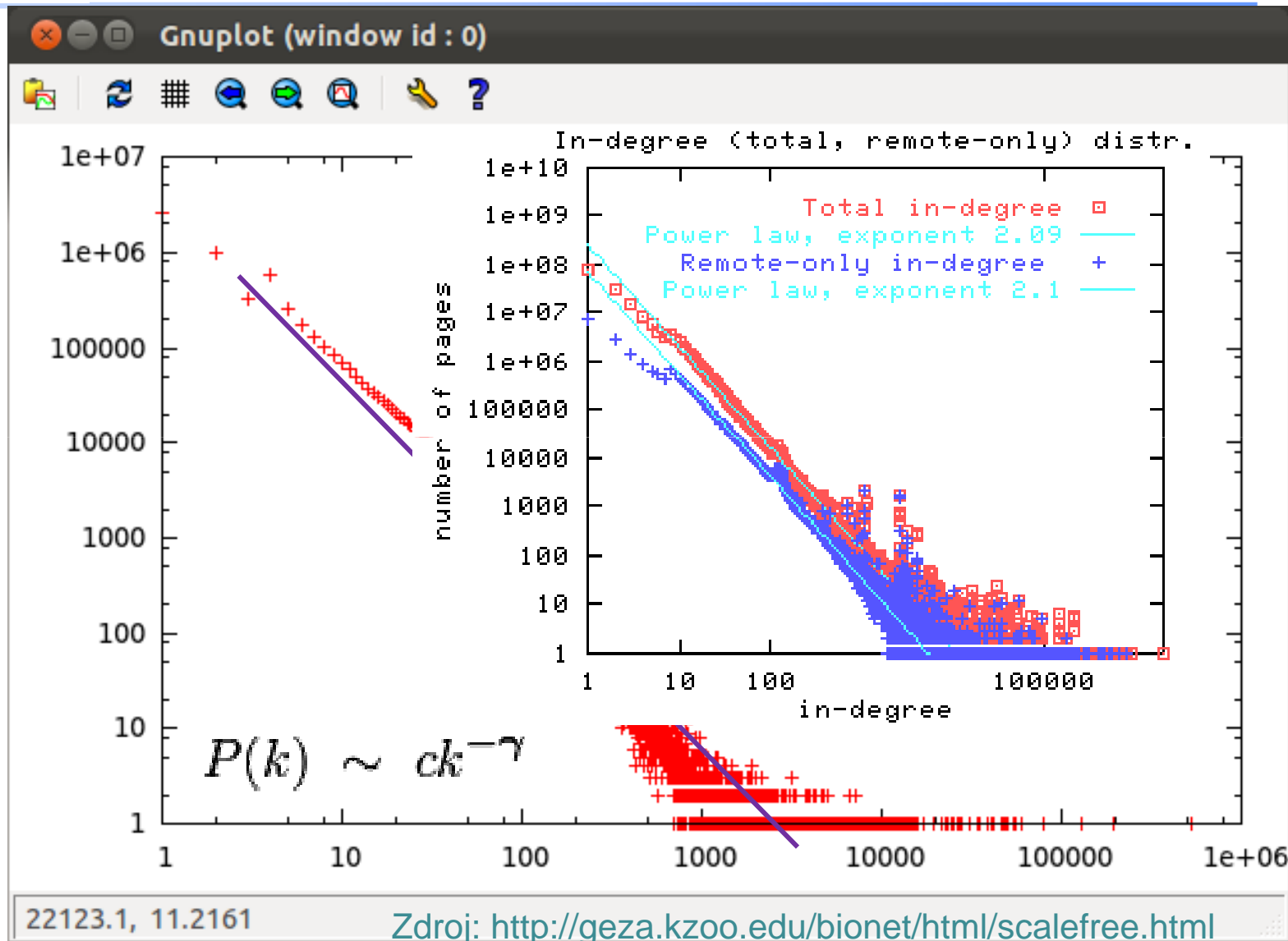




Európska únia
Európsky fond regionálneho rozvoja

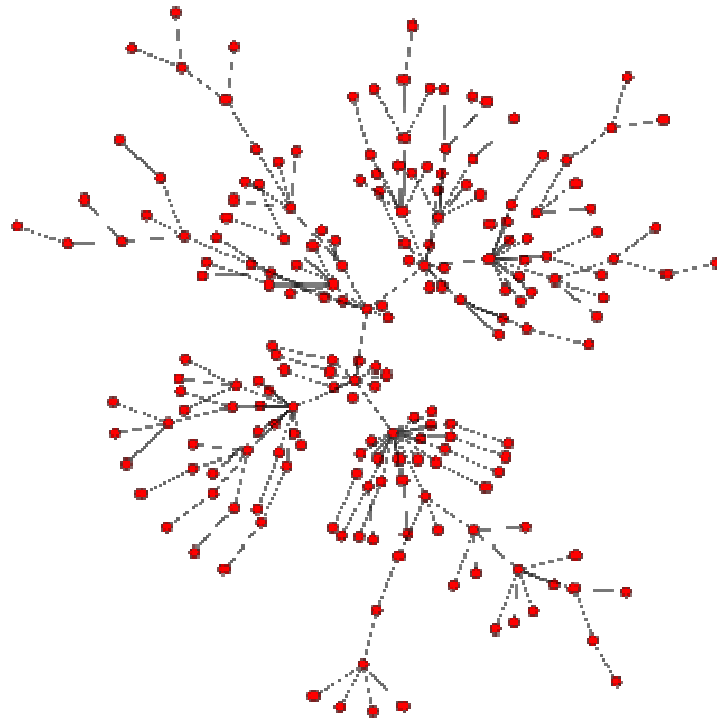
Distribúcia stupňov v reálnych sieťach - log škála (príklad DBpedia)



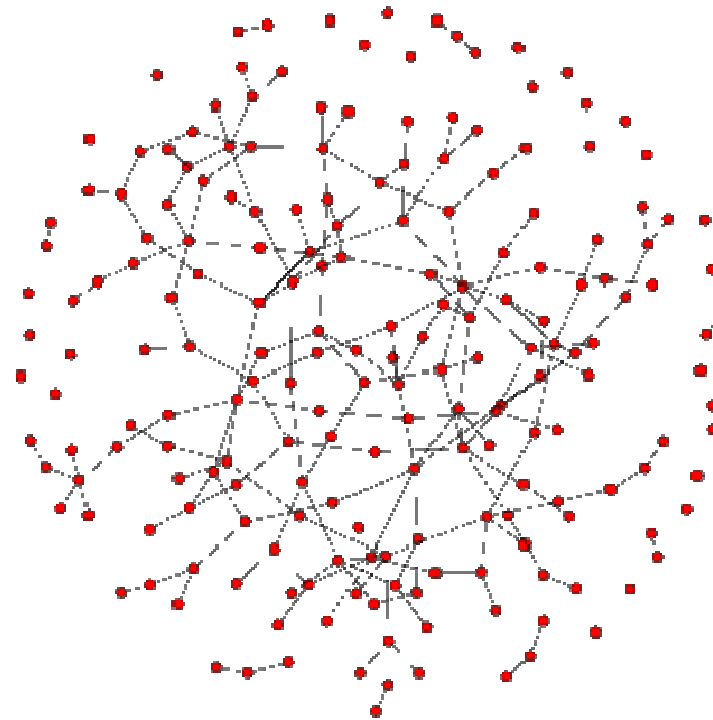




Sieť s mocninovou dist. stuňov



Sieť s binomickou dist. stuňov

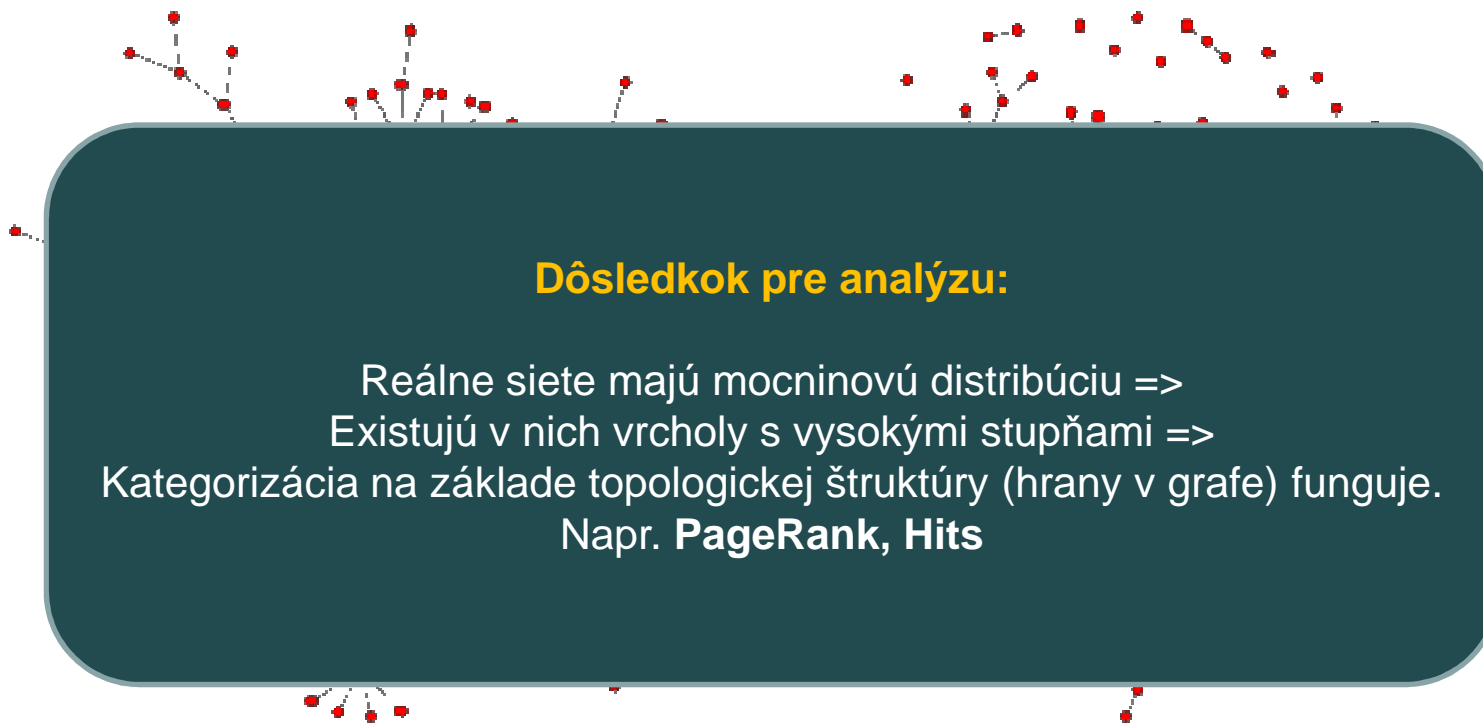


Zdroj: <http://geza.kzoo.edu/bionet/html/scalefree.html>



Sieť s mocninovou dist. stuňov

Sieť s binomickou dist. stuňov



Zdroj: <http://geza.kzoo.edu/bionet/html/scalefree.html>

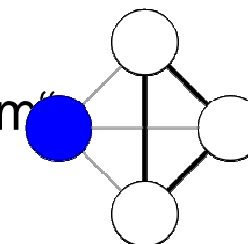


- Objavené štúdiom sociálnych sietí v sociológii
- Milgramov experiment
 - (60-te roky; úloha: doručenie listov cez sociálnu sieť od náhodných ľudí v Nebraske prímateľom v Chicagu; priemer 6 krokov pre)
- Objav: sociálne siete majú *krátku* dĺžku najkratších ciest medzi náhodne vybranými
- Potvrdené aj skúmaním počítačových sociálnych sietí
 - **Analýza MSM site – instant messaging**
 - **Priemerná dĺžka najkrajšej siete medzi náhodne vybranými uzali: 6,6**
 - *[Jure Leskovec, Eric Horvitz: Planetary-scale views on a large instant-messaging network. WWW 2008]*
- Priemer siete (aj priemerná dĺžka cesty medzi 2 vrcholmi) sa znižuje pri zväčšovaní siete
 - V protiklade k náhodným sieťam

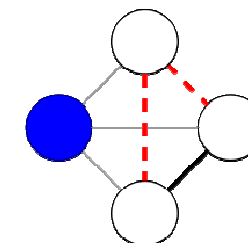


- Siete malého sveta často obsahujú kliky, alebo „skoro kliky“
- Efekt „moji priatelia v sociálnej sieti sú často priatelia navzájom“
- Matematicky to možno zachytiť pomocou zhlukovacieho koeficientu
- Lokálny zhlukovací koeficient:

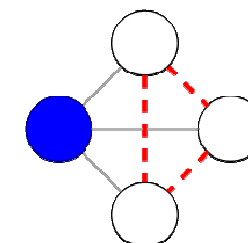
$$C_i = \frac{|\{e_{jk}\}|}{k_i(k_i - 1)} : v_j, v_k \in N_i, e_{jk} \in E.$$



$$c = 1$$



$$c = 1/3$$



$$c = 0$$

Zdroj: http://en.wikipedia.org/wiki/Clustering_coefficient



- Mocninová distribúcia stupňov vrcholov
- Malá vzdialenosť medzi náhodnými uzlami v sieti (malý priemer grafu)
- Priemer grafu sa znižuje pri rozširovaní siete
- Vzor zhukovania v sieti: vysoký zhukovací koeficient
- Dôsledky:
 - Algoritmy na spracovanie/dolovanie grafov fungujú vďaka týmto vlastnostiam
 - Rozličné typy sietí z reálneho sveta majú podobné vlastnosti
 - Algoritmy navrhnuté pre jeden konkrétny typ sietí (napr. soc. siete) budú pravdepodobne dávať zmysluplné výsledky aj na iných sieťach s podobnými matematickými vlastnosťami
 - Mocninová distribúcia stupňov – pri traversovaní grafu do šírky už pri nízkom počte hopov je nutné prejsť značnú časť siete



- Rekurzívne počítanie mier centralít:
 - Odhadnúť dôležitosť vrchlov v topológii siete
 - PageRank
 - HITS
- Detekcia komunít
 - Identifikovať skupiny vrcholov silne prepojené medzi sebou a slabšie prepojené s ostatnými komunitami
 - Prekrývajúce sa komunity
- Propagácia v sieťach
 - Šírenie informácií v sieťach
 - Kaskádové správanie, propagácia vírusov
- Klasifikácia objektov na základe liniek
- Predikcia vzniku liniek
- Objavovanie častých vzorov

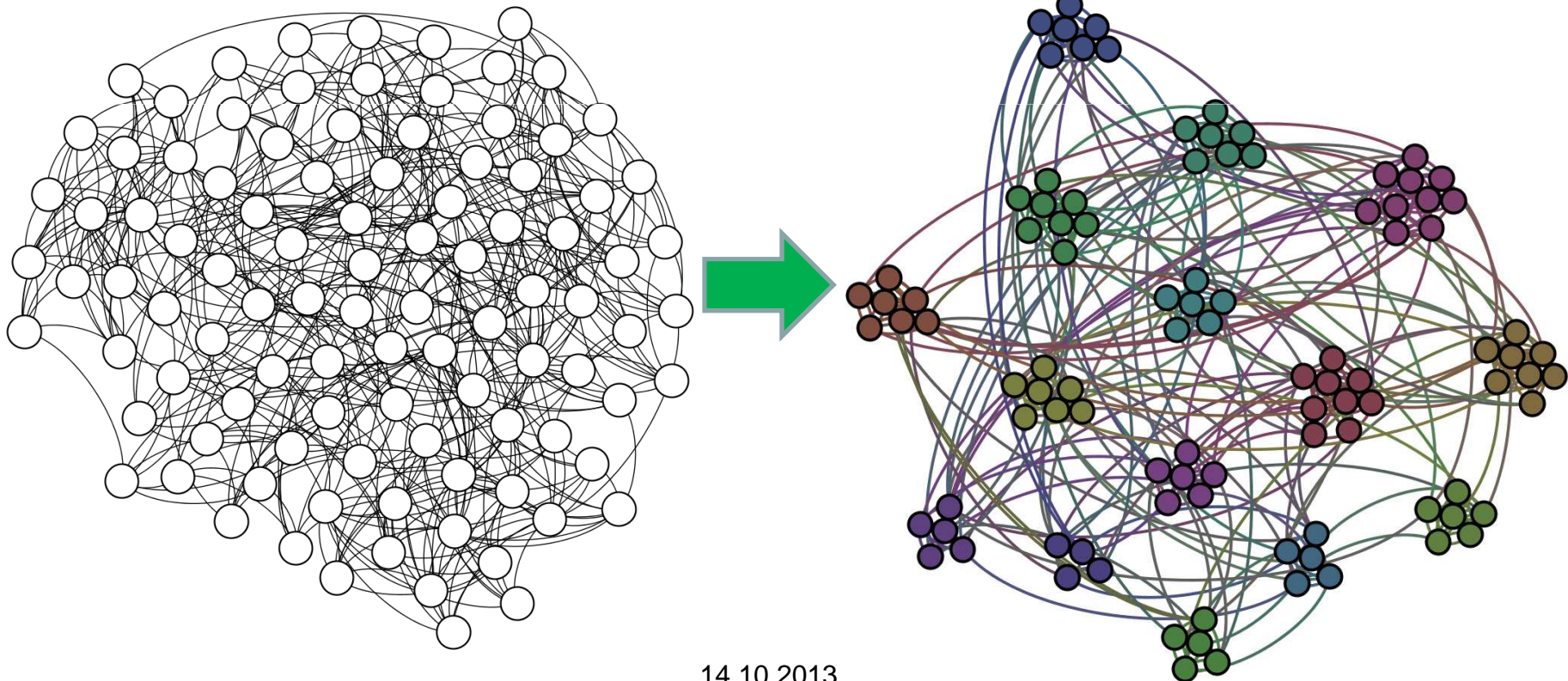


Európska únia
Európsky fond regionálneho rozvoja

Detekcia komunít



- Identifikácie zhukov uzlov siete silne prepojených medzi sebou a menej silne prepojených so zvyškom siete
- Detekcia komunít v grafoch má za cieľ identifikovať moduly a ich prípadnú hierarchickú organizáciu



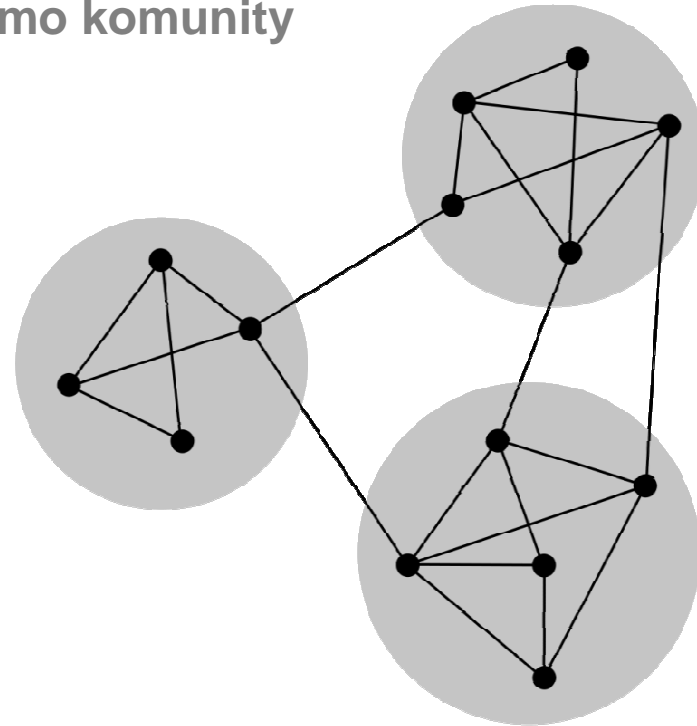
14.10.2013



- Neexistuje konsenzus na formálnej definícii :
 - problému detekcie komunit
 - komunitnej štruktúry grafu (community structure)
 - komunity ako takej
- Najčastejšie používané definície komunity / zhľuku grafu sú založené na počte hrán v rámci komunity (hustota) a počte hrán vedúcich mimo komunitu
- Definície komunitnej štruktúry sú často silne závislé na problémovej doméne a vlastnostiach analyzovaných grafov.



- Najčastejšie používaná definícia:
 - Komunita je zhluk uzlov, kde počet vnútorných hrán v komunite je väčší ako počet vnokajších hrán – mimo komunity



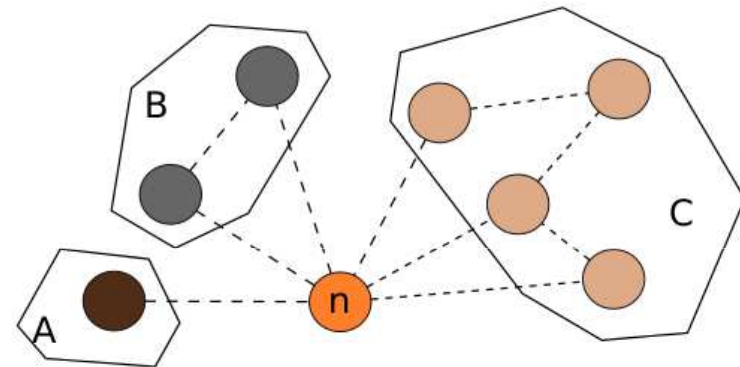
[M.E.J. Newman. The structure and function of complex networks. SIAM REVIEW 45, 2003.]



- Alternatívna definícia:

- Komunita je zhluk uzlov, kde pre každý uzol platí, že počet hrán uzla v rámci komunity je väčší ako počet hrán uzla smerujúcich do inej komunity.
- $G = (V, E)$

$$aff(n, C) = \sum_{i \in C} w(e_{n,i}) : e_{n,i} \in E$$



$$\bigcup_{k \in \langle 1, \dots, m \rangle} C_k = V \quad \text{and} \quad \bigcap_{k \in \langle 1, \dots, m \rangle} C_k = \emptyset$$

$$\forall j \in C_k, aff(j, C_k) \geq \max\{aff(j, C_l), C_l \in \gamma\}$$

- Y. Hu, H. Chen, P. Zhang, M. Li, Z. Di, and Y. Fan. Comparative definition of community and corresponding identifying algorithm. *Phys. Rev. E*, 78(2):026121, Aug 2008
- Marek Ciglan, Kjetil Nørvåg: Fast detection of size-constrained communities in large networks, proceedings of WISE'10, LNCS Volume 6488/2010



- Dôsledky neexistencie konsenzu na formálnej definícií problému:
 - Rôzni autori používajú rôzne definície komunit
 - Typické riešenie problému, 2 kroky:
 - Špecifikácia metriky kvality komunitnej štruktúry
 - Algoritmická technika pre zhlukovanie grafu optimalizujúca danú metriku
- Veľké množstvo publikácií
 - Prehľadová práca [Fortunato10]
 - 457 referencií
 - Citovaná viac ako 1800 krát (podľa služby google scholar)
 - Publikačný pretlak - Počet publikovaných algoritmov na detekciu komunit možno rátať v stovkách
 - Dopyt „community detection“ v google scholar - 8000+ výsledkov
 - Dopyt „community detection algorithm“ v google scholar - 1500+ výsledkov

[Fortunato10] S. Fortunato (2010). "Community detection in graphs". Phys. Rep. 486 (3-5): 75–174.
doi:10.1016/j.physrep.2009.11.002

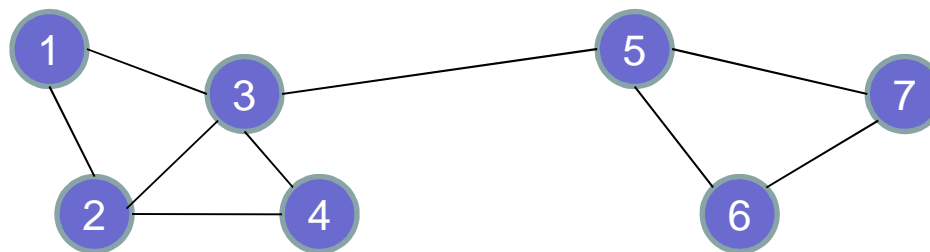


- Algoritmy založené na delení siete [Girvan, Newman]
- Algoritmy založené modularite
 - Greedy aproximácie
 - Simulované žíhanie
 - Extremal optimization
- Algoritmy založené na spektrálnej analýze grafu
 - Využíva spektrum grafu (vlastné čísla a vektory maticovej reprezentácie grafu)
- Dynamické algoritmy
 - Náhodné pochôdzky (random walk)
 - Spinové modely
- Modely založené na štatistickej inferencii



- Založený na koncepte centrality spojitosti hrán (betweenness centrality)

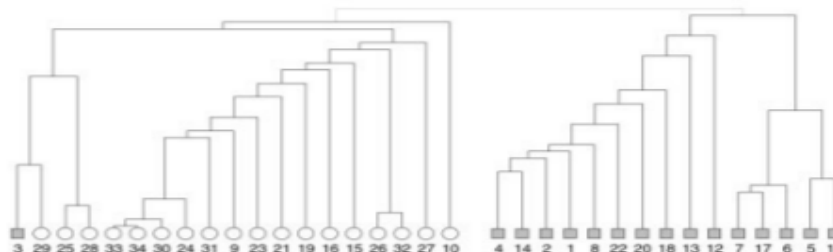
- Centralita spojitosti pre e :
 - Počet najkratších ciest prechádzajúcich hranou e



- Algoritmus:

1. Vypočítaj centralitu spojitosti pre všetky hrany grafu
2. Hrana s najvyššou mierou centrality je odstránená
3. Prepočítaj mieru centrality pre hrany ovplynené odobratím hrany
4. Krok 2 a 3 je opakovaný do odobratia poslednej hrany

- Výsledkom algoritmu je dendrogram



Girvan M. and Newman M. E. J., Community structure in social and biological networks, Proc. Natl. Acad. Sci. USA 99, 7821–7826 (2002)



- Typický návrh algoritmu, 2 kroky:
 - Špecifikácia *metriky* kvality komunitnej štruktúry
 - Algoritmická technika pre zhlukovanie grafu optimalizujúca danú metriku
- Metriky zamerané na:
 - Početnosť hrán v rámci zhluku
 - Početnosť hrán vedúcich mimo zhluku
 - Kombinácia oboch

• Notácia: $G = (V, E)$

$$C : C \subseteq V$$

$$n_c = |C|$$

$$m_c = |F| : F \subseteq E \wedge \forall (i, j) \in F : i, j \in C$$

$$o_c = |G| : g \subseteq E \wedge \forall (i, j) \in F : i \in C \wedge j \notin C$$

$$d_v = |K| : K \subseteq E \wedge \forall (i, j) \in K : i = v \wedge j \in V$$



- Vnútorná hustota (Internal density)

- Vnútorná hustota komunity C

$$f(C) = \frac{m_c}{(n_c \times (n_c - 1)) / 2}$$

- Priemerný stupeň (Average degree)

- Priemerný počet vnútorných hrán na uzol

$$f(C) = \frac{2 \times m_c}{n_c}$$

- Priemerný počet trojuholníkov

- Priemerný počet vnútorných trojuholníkov na uzol

$$f(C) = \frac{|\{(v:v \in C, \{u,w \in C \wedge (v,u) \in E \wedge (v,w) \in E \wedge (u,w) \in E\} \neq \emptyset)\}|}{n_c}$$



- Vnútorná hustota (Internal density)

- Vnútorná hustota komunity C

$$f(C) = \frac{m_c}{(n_c \times (n_c - 1)) / 2}$$

[F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi. Defining and identifying communities in networks. PNAS, 101(9), 2004.]

- Priemerný stupeň (Average degree)

- Priemerný počet vnútorných hrán na uzol

$$f(C) = \frac{2 \times m_c}{n_c}$$

[J. Yang and J. Leskovec. Defining and Evaluating Network Communities based on Ground-Truth. In: ICDM, 2012.]

- Priemerný počet trojuholníkov

- Priemerný počet vnútorných trojuholníkov na uzol

$$f(C) = \frac{|\{(v:v \in C, \{u, w \in C \wedge (v, u) \in E \wedge (v, w) \in E \wedge (u, w) \in E\} \neq \emptyset)\}|}{n_c}$$



- Expanzia (Expansion)
 - Počet hrán mimo komunity na uzol

$$f(C) = \frac{o_c}{n_c}$$

- Pomer rezu (cut ratio)
 - Pomer existujúcich hrán idúcich mimo komunity k počtu možných hrán

$$f(C) = \frac{o_c}{n_c \times (n \times n_c)}$$



- Expanzia (Expansion)
 - Počet hrán mimo komunity na uzol

$$f(C) = \frac{o_c}{n_c}$$

[F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi. Defining and identifying communities in networks. PNAS, 101(9), 2004.]

- Pomer rezu (Cut ratio)
 - Pomer existujúcich hrán idúcich mimo komunity k počtu možných hrán

$$f(C) = \frac{o_c}{n_c \times (n \times n_c)}$$

[Fortunato10] S. Fortunato (2010). "Community detection in graphs". Phys. Rep. 486 (3-5): 75–174. doi:10.1016/j.physrep.2009.11.002]



- Vodivosť (Conductance)
 - Pomer počtu hrán vedúcich mimo komunitu ku celkovému počtu hrán

$$f(C) = \frac{o_c}{2m_c + o_2}$$

- Maximálny vonkajší stupeň (Maximum out degree fraction)
 - Maximálny pomer von idúcich hrán z vrchola komunity k jeho stupňu

$$f(C) = \max_{v \in C} \frac{|\{(v, u) \in E \wedge u \notin C\}|}{d_v}$$

- Priemerný vonkajší stupeň (Average out degree fraction)
 - Priemerný pomer von idúcich hrán z vrchola komunity k jeho stupňu

$$f(C) = \frac{1}{n_c} \sum_{v \in C} \frac{|\{(v, u) \in E \wedge u \notin C\}|}{d_v}$$



- Vodivosť (Conductance)
 - Pomer počtu hrán vedúcich m

$$f(C) = \frac{o_c}{2m_c + o_2}$$

Fan Chung. Spectral Graph Theory. CBMS Lecture Notes 92, AMS Publications, 1997.

án

- Maximálny vonkajší stupeň (Maximum out degree fraction)
 - Maximálny pomer von idúcich hrán z vrchola komunity k jeho stupňu

$$f(C) = \max_{v \in C} |\{(v, u) \in E \wedge u \notin C\}|$$

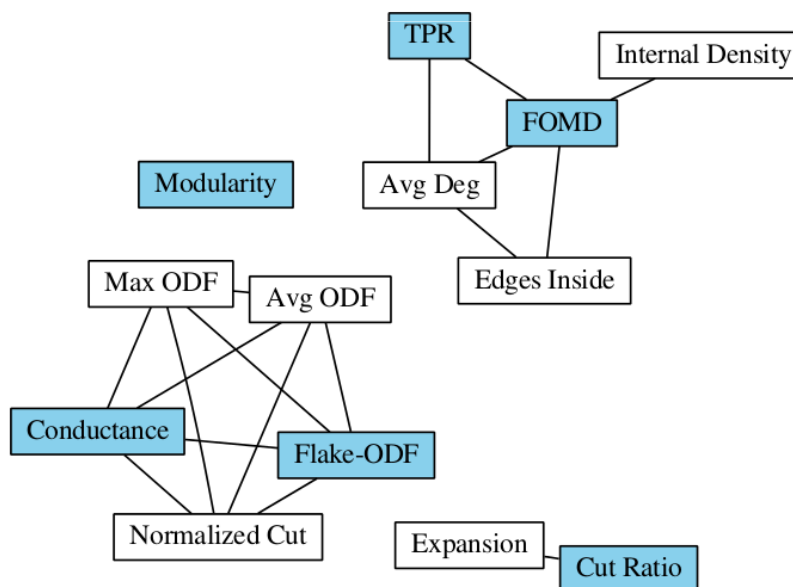
G. Flake, S. Lawrence, and C. Giles. Efficient identification of web communities. In KDD '00, pp 150–160, 2000

- Priemerný vonkajší stupeň (Average out degree fraction)
 - Priemerný pomer von idúcich hrán z vrchola komunity k jeho stupňu

$$f(C) = \frac{1}{n_c} \sum_{v \in C} \frac{|\{(v, u) \in E \wedge u \notin C\}|}{d_v}$$



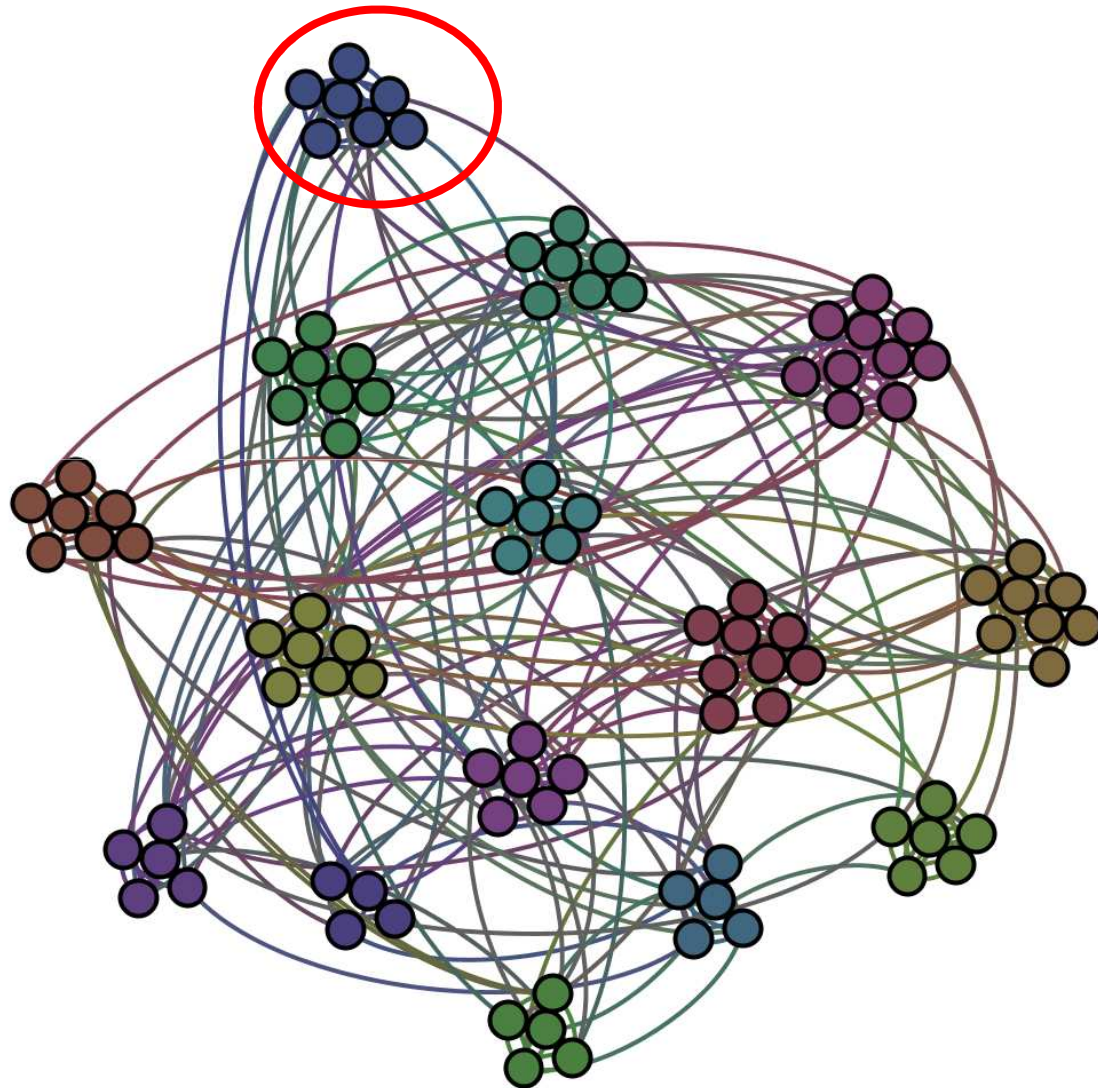
- Aký je vzťah medzi rozličnými metrikami?
- Zdroj: [J. Yang and J. Leskovec. *Defining and Evaluating Network Communities based on Ground-Truth. In: ICDM, 2012*]
- Autori použili reálne siete s explicitne definovanými komunitami
- Pre komunity vypočítali hodnoty jednotlivých metrick
- Porovnali koreláciu





Európska únia
Európsky fond regionálneho rozvoja

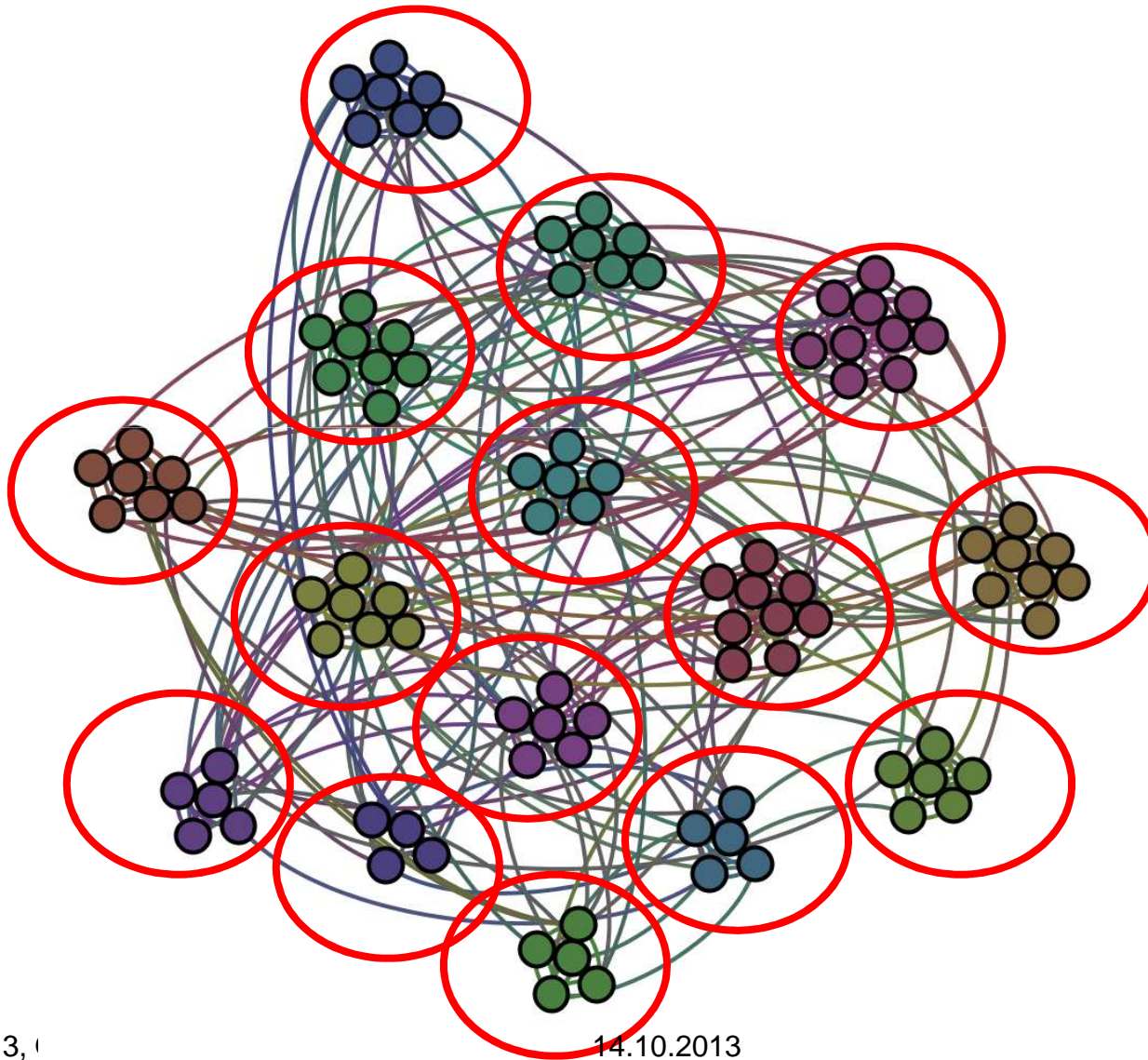
Metriky – globálne metriky rozdelenia siete





Európska únia
Európsky fond regionálneho rozvoja

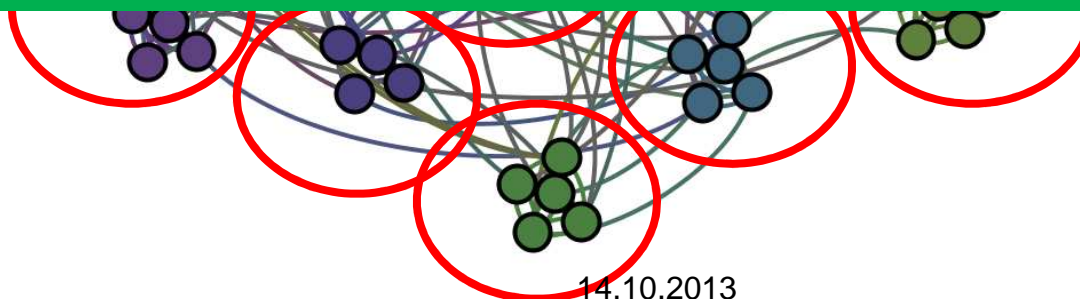
Metriky – globálne metriky rozdelenia siete





Možno použiť lokálne metriky kvality komunity na posúdenie / porovnanie globálnej komunitnej štruktúry celej siete?

- **Problematické sú hraničné prípady:**
 - pre vnútornú konektivitu: rozdelenie do skupín po dvoch uzloch s jednou hranou maximalizuje takúto metriku
 - pre vonkajšiu konektivitu: rozdelenie do skupín na základe spojených komponentov maximalizuje takúto metriku





- Zdoj: *[M.E.J. Newman and M. Girvan. Finding and evaluating community structure in networks. Physical Review E 69(02), 2004.]*
- Modularita – metrika ohodnocujúca vhodnosť rozdelenia celej siete
- Myšlienka:
 - Porovnať počet hrán ležiacich v zhlukoch s očakávaným počtom hrán v referenčnom modeli siete s rovnakým počtom vrcholov a hán
 - Referenčný model: obvykle náhodná sieť
- Modularita
 - $Q = (\text{početnosť hrán vo vnútri komunit}) - (\text{očakávaná početnosť hrán v rovnakom rozdelení do skupín v referenčnom modeli (v náhodnej sieti)})$



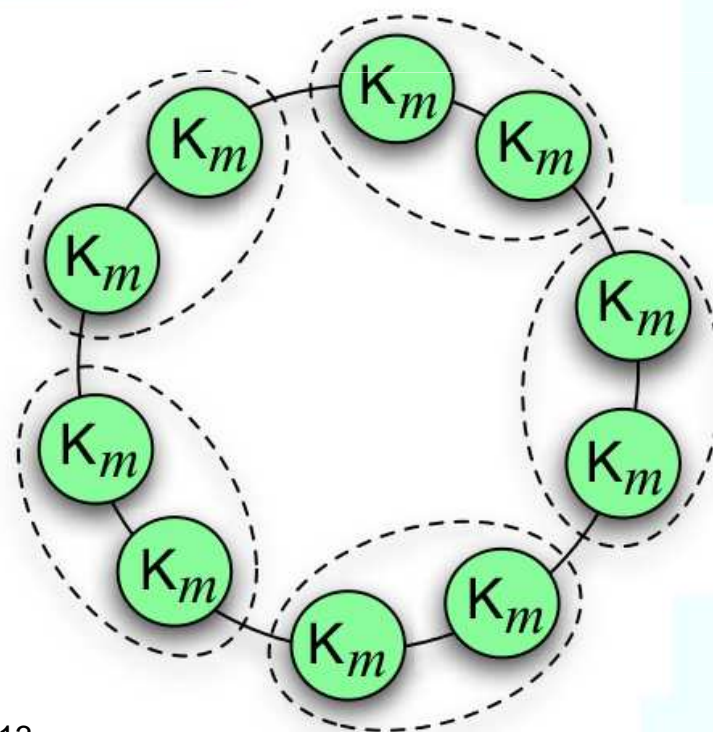
- Modularita

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

- $A_{i,j}$ – incidenčná matica
- k_i, k_j – stupne vrcholov i a j
- m – počet hrán grafu
- $\Delta(c_i, c_j) = 1$ iff i a j ležia v rovnakom zhluku, ináč 0
- Modularita vhodná ako
 - Metrika kvality rozdelenia
 - Metrika na porovnanie rôznych rozdelení, algoritmov



- Práca: S. Fortunato and M. Barthelemy. Resolution limit in community detection. PNAS 104(1), 2007
- Problém modularity: **limit rozlíšenia**
- Optimalizácia modularity nerozozná malé komunity (relatívne, na počet hrán siete)
- K_m – klik o m hranách
- $m < \sqrt{|E|}$
 - Maximálne skóre modularity odpovedá rozdeleniu s viacerými klikmi v komunitách





- Práca: *[B. H. Good, Y.-A. de Montjoye and A. Clauset: The performance of modularity maximization in practical contexts. Physical Review E 81, 046106 (2010)]*
- Modularita vykazuje extrémnu degeneráciu - Exponenciálne množstvo rôznych rozdelení grafu do komunít nadobúda vysoké hodnoty modularity



- Neexistencia konsekvencie na formálnej špecifikácii problému
 - Rôzne definície / intuitívne, neformálne definície
 - Veľké množstvo navrhnutých algoritmov
- Metriky
 - Lokálne – množstvo rôznych definícií
 - Lokálne metriky – zle transformovateľné na globálnu úroveň rozdelenia celej siete
- Modularita
 - Metrika ohodnotenia rozdelenia celej siete
 - Limit rozlíšenia – relatívne malé komunity nie sú odhalené
 - Problém degenerácie – veľa rozdelení s vysokým skóre modularity



- Alternatívny prístup k vyhodnocovaniu presnosti algoritmov na detekciu komunit
- Majme grafy so známou štruktúrou komunit
- Výsledky zhlukovacieho algoritmu porovnáme s referenčným rozdelením
- Grafy možno generovať

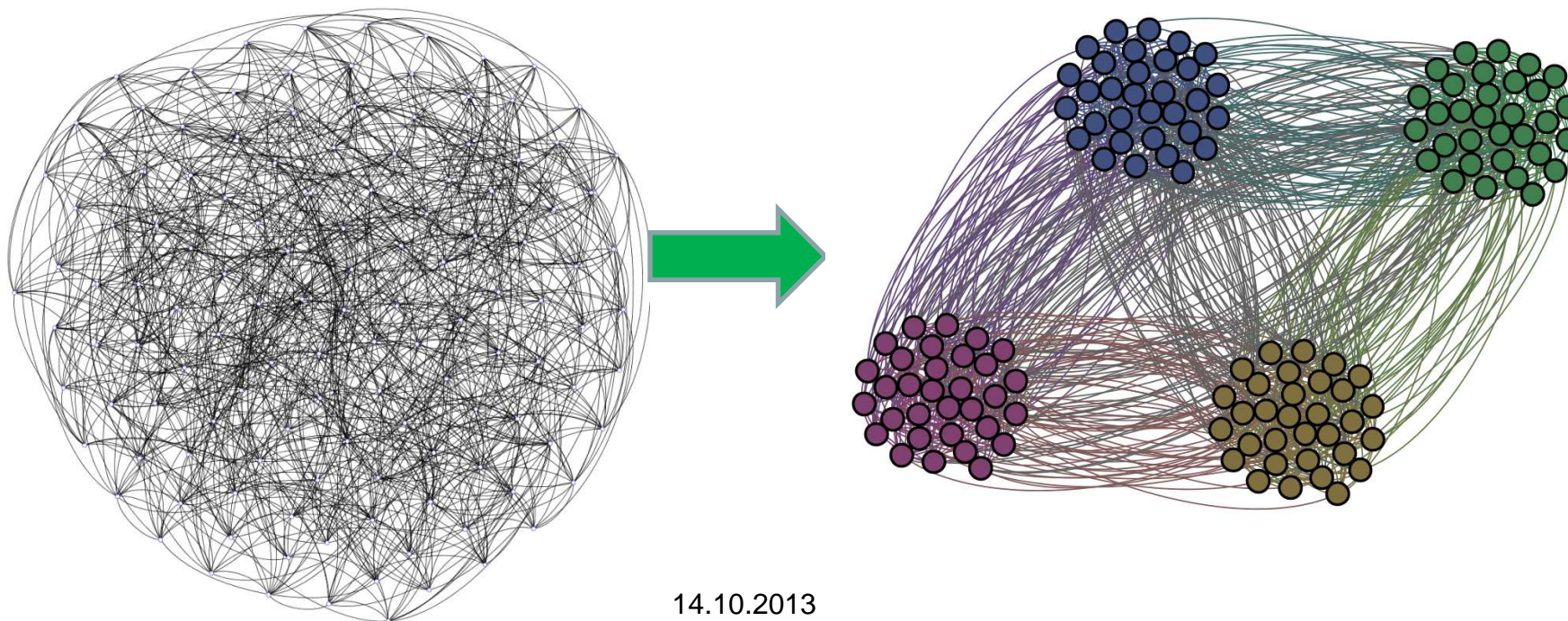


Európska únia
Európsky fond regionálneho rozvoja

Girvan-Newman benchmark grafy



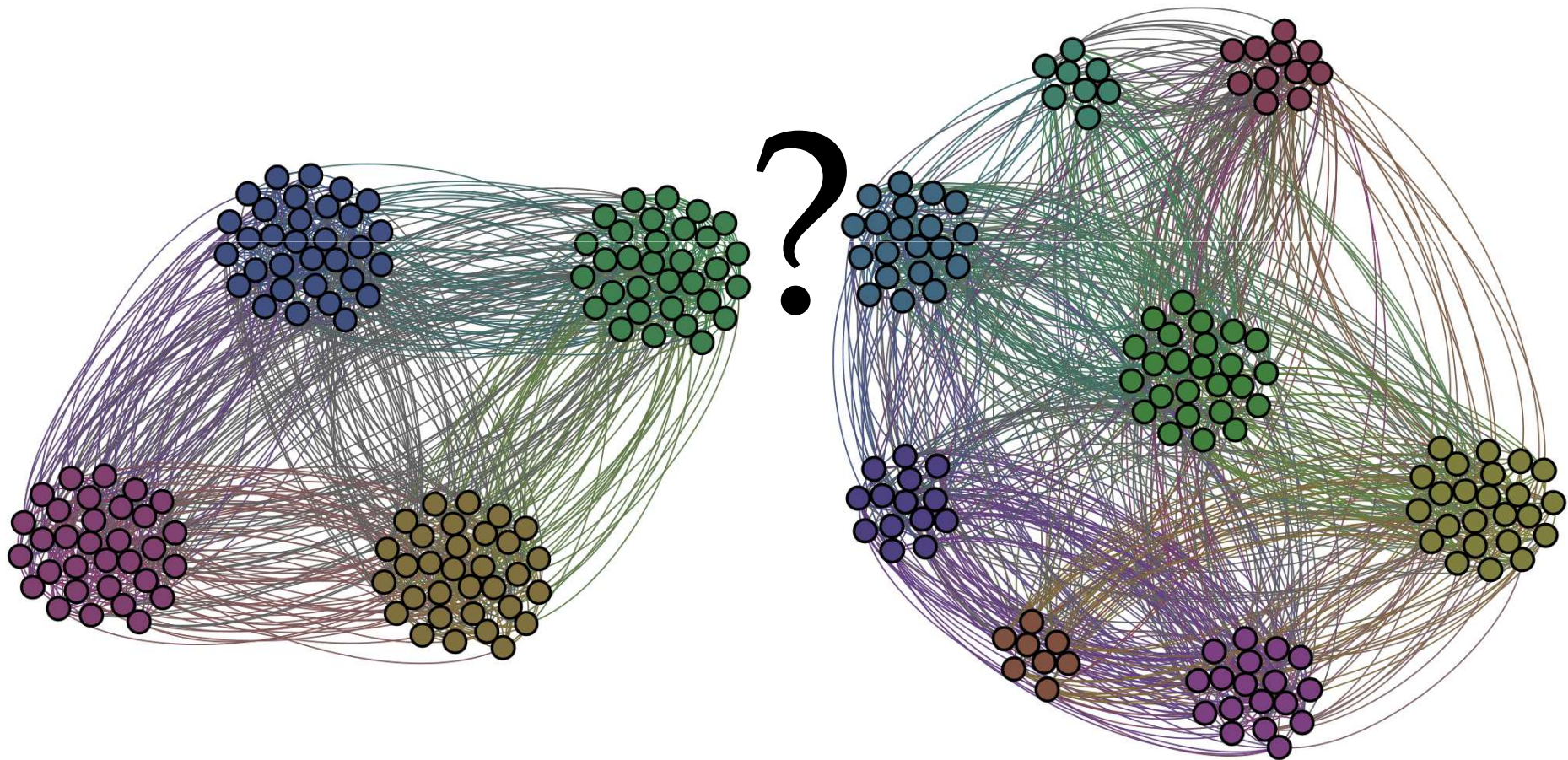
- Práca: [M. Girvan and M. E. J. Newman, *Community structure in social and biological networks*, *Proc. Natl. Acad. Sci.* 99, 7821 (2002).]
- Benchmark grafy
 - 128 vrcholov
 - Rozdelných do 4-och komúní
 - Každá komunita: 32 uzlov
 - Priemerný stupeň uzla: 16





Európska únia
Európsky fond regionálneho rozvoja

Porovnanie detekovaného a referenčného rozdelenia





- Práca: [L. Danon, J. Duch, A. Diaz-Guilera, and A. Arenas. *Comparing community structure identification. Journal of Statistical Mechanics: Theory and Experiment, Oct 2005*]
- Navrhnutá miera „Normalized mutual information“ (information theory)
- Výpočet založený na prienikovej matici N (confusion matrix)
 - Riadky reprezentujú referenčné komunity
 - Stĺpce reprezentujú detekované komunity
 - Prvky matice: N_{ij} = Počet prvkov referenčnej komunity i , ktoré sa nachádzajú v detekovanej komunite j

$$I(A, B) = \frac{-2 \sum_{i=1}^{c_A} \sum_{j=1}^{c_B} N_{ij} \log \left(\frac{N_{ij} N}{N_{i.} N_{.j}} \right)}{\sum_{i=1}^{c_A} N_{i.} \log \left(\frac{N_{i.}}{N} \right) + \sum_{j=1}^{c_B} N_{.j} \log \left(\frac{N_{.j}}{N} \right)}$$



- Jaccardova podobnosť:
 - štatistika určujúca podobnosť dvoch množín

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- Porovnáva dve množiny
- Použitie na porovnanie dvoch rozdelení do komunit:
 1. Pre každú referenčnú komunitu
 2. nájdí detekovanú komunitu s najväčším prienikom
 3. Vypočítaj Jaccardov index
 4. Výsledok: priemerný jaccardov index / alebo medián

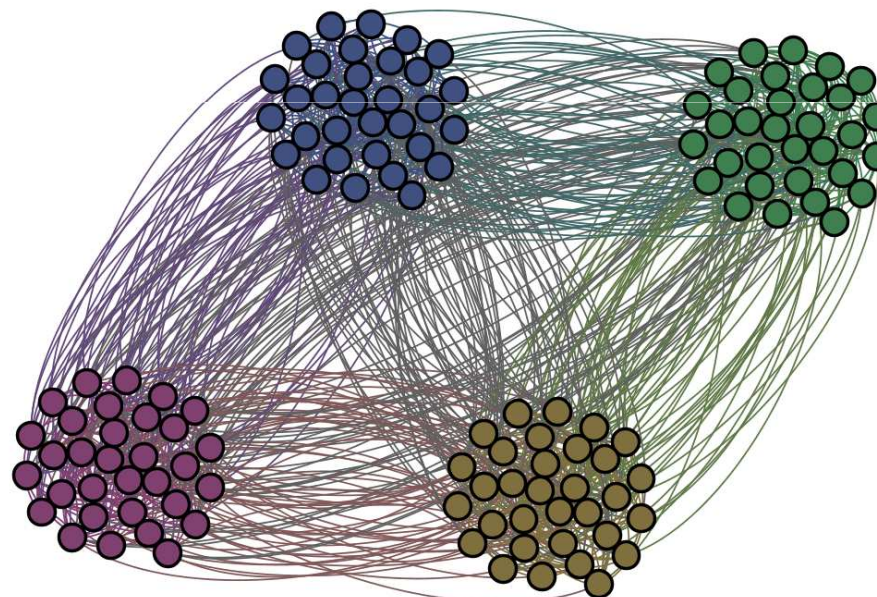


Európska únia
Európsky fond regionálneho rozvoja

GN benchmark nets - kritika

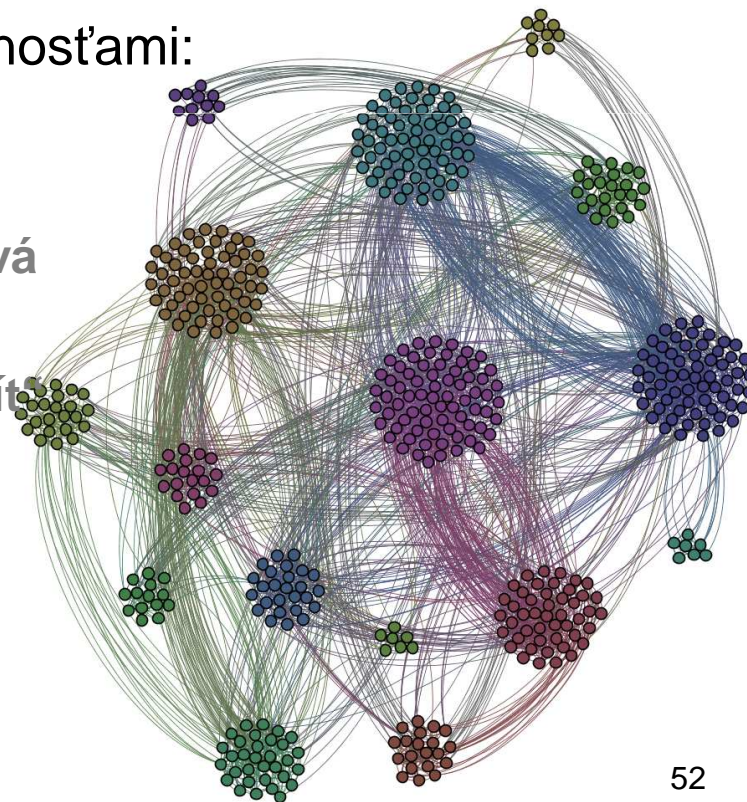


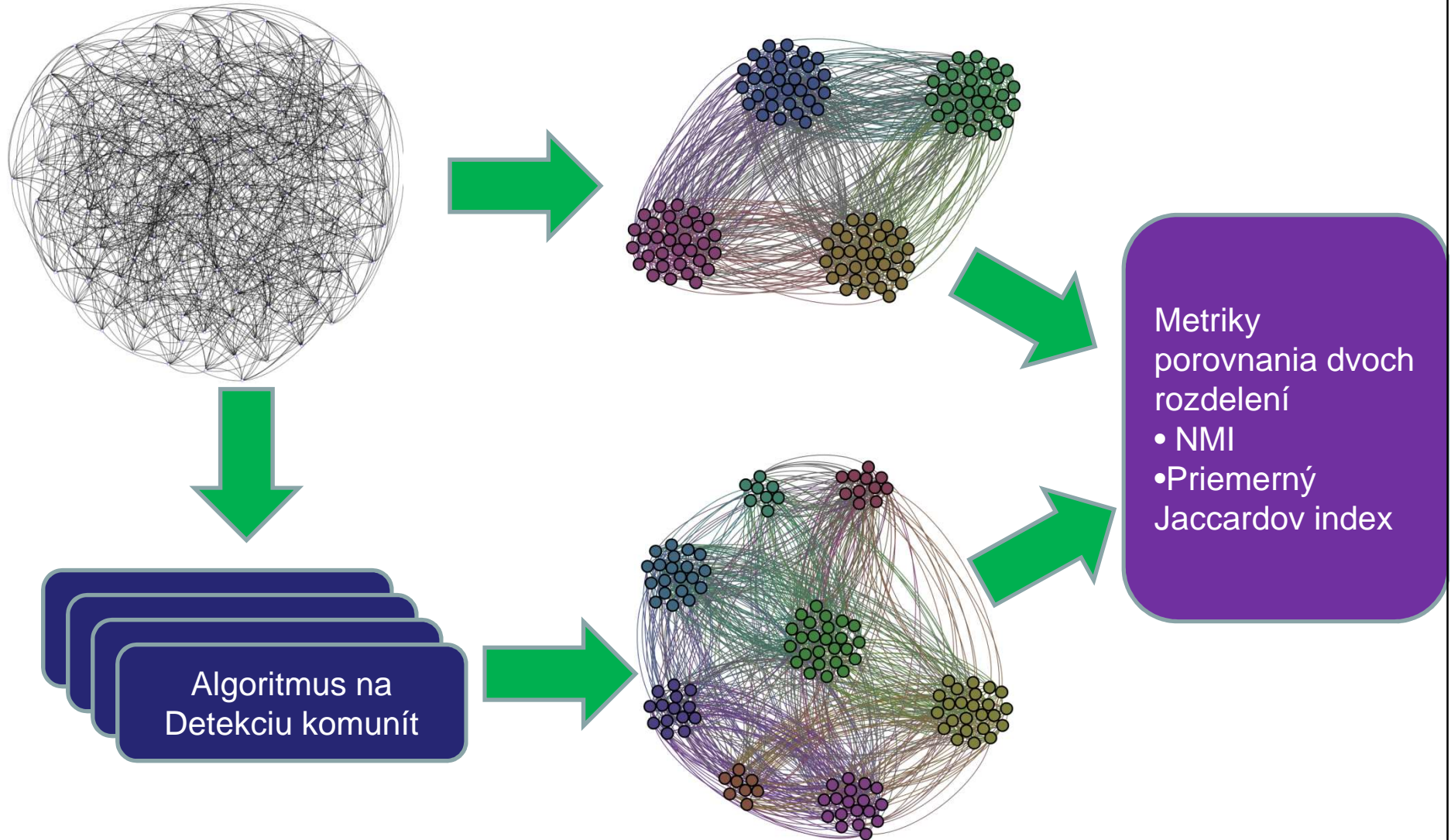
- GN benchmark siete:
 - Nerealistické
 - Malé – 128 uzlov
 - Malý počet komunit
 - Rovnaká veľkosť komunit
 - Nerealistická distribúcia uzlov





- Práca: [A. Lancichinetti and S. Fortunato. Benchmarks for testing community detection algorithms on directed and weighted graphs with overlapping communities. *Phys. Rev. E*, 80(1), 2009.]
- Snaha odstániť nedostatky GN-benchmark grafov
- Generovanie sietí s viac realistickými vlastnosťami:
 - Založené na planted l-partition model
 - Distribúcia stupňov – mocninová
 - Distribúcia veľkostí komunít – mocninová
 - Definovateľný zhukovací koeficient
 - Parametrizovateľná „viditeľnosť komunít“







Európska únia
Európsky fond regionálneho rozvoja

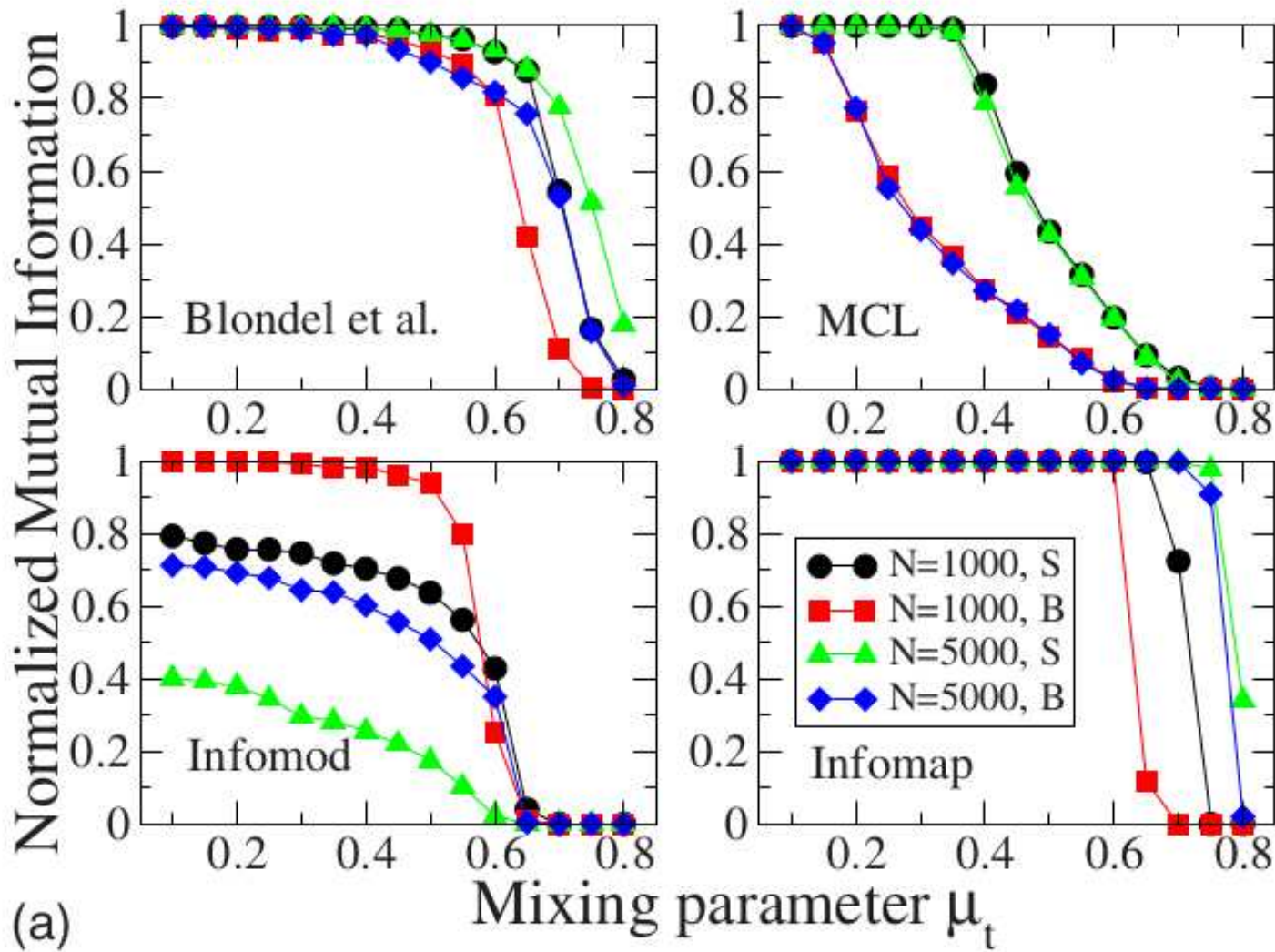
Modifikácie základnej úlohy



- Detekcia prekrývajúcich sa komunít
 - Jeden vrchol patrí do viacerých komunít
- Siete s atribútmi
 - Vrcholy a/alebo hrany majú atribúty
 - Dodatočná informácia k topológii grafu
- Detekcia komunity / komunít pre daný vrchol
- Hierarchická detekcia komunít

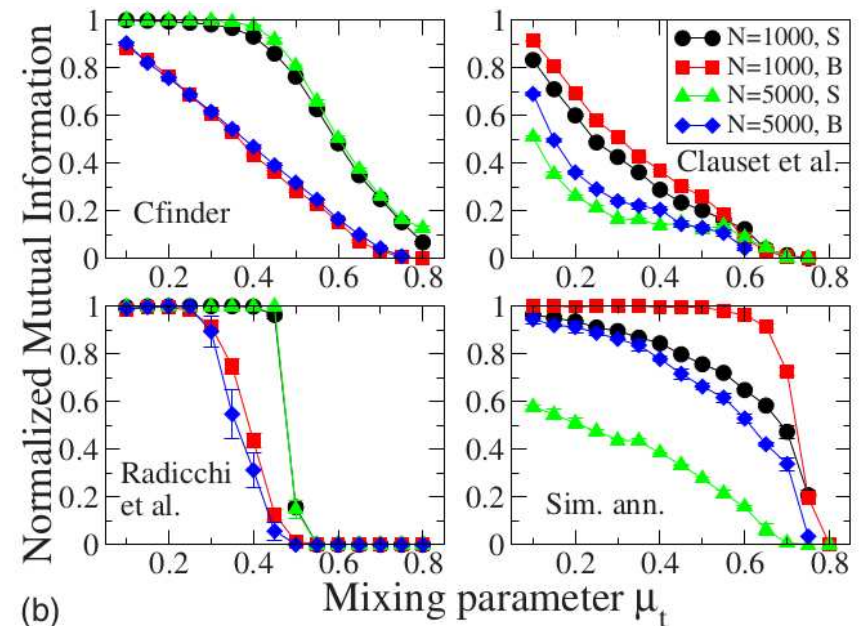
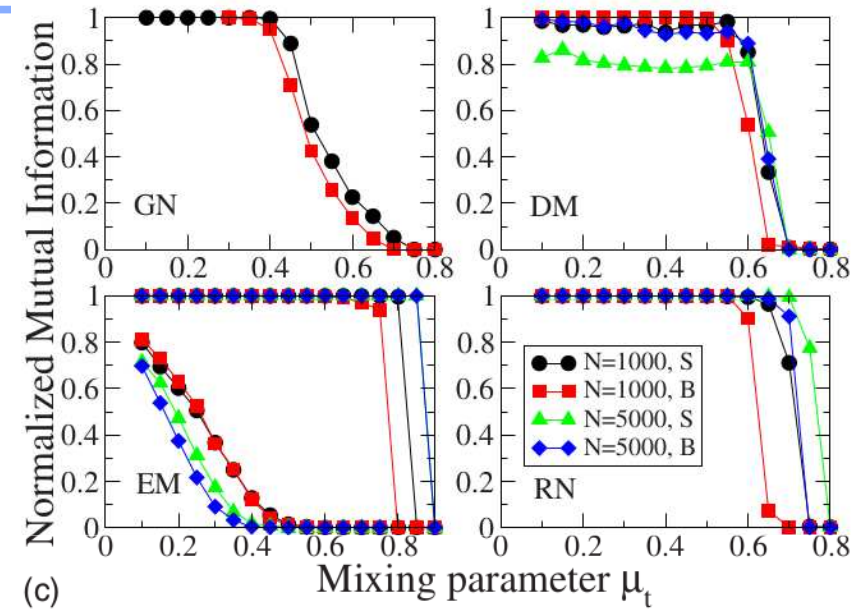
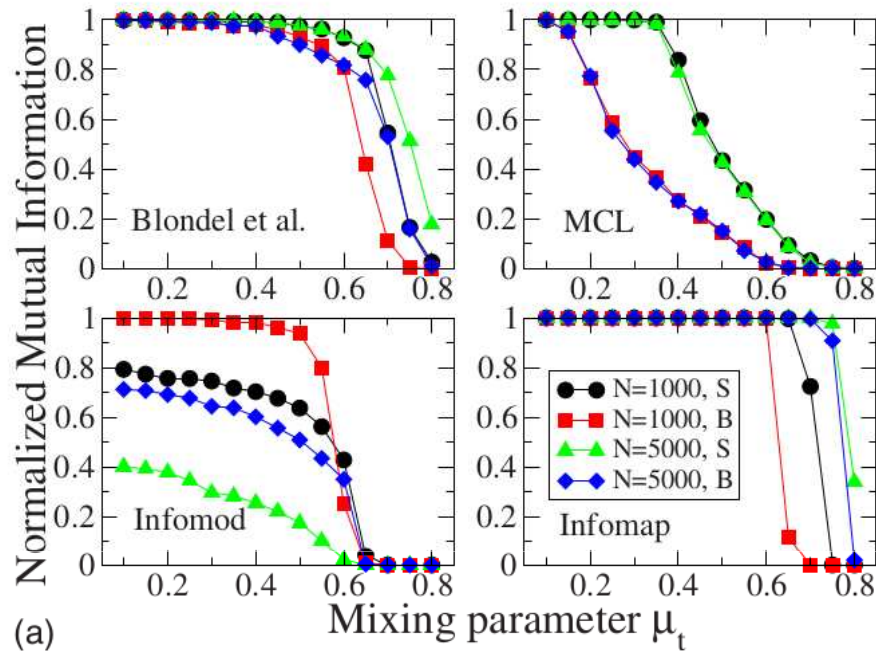


- Práca: [A. Lancichinetti and S. Fortunato. *Community detection algorithms: A comparative analysis*. *Phys. Rev. E*, 80(5):056117, Nov 2009.]
- Porovnanie populárnych algoritmov na detekciu komunit
- Použitie LFR bechmark grafov
- Porovnanie na rôznych nastaveniach generovaných sietí



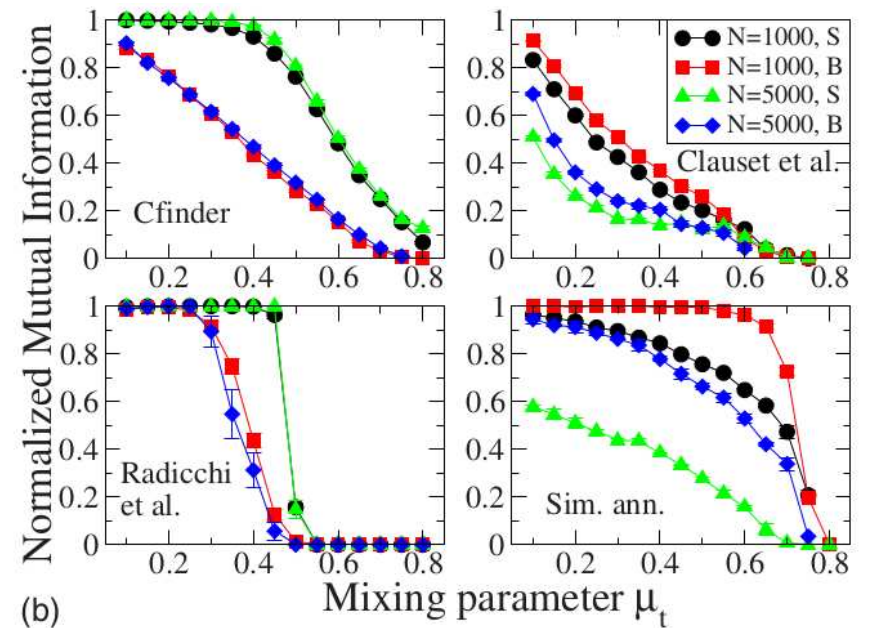
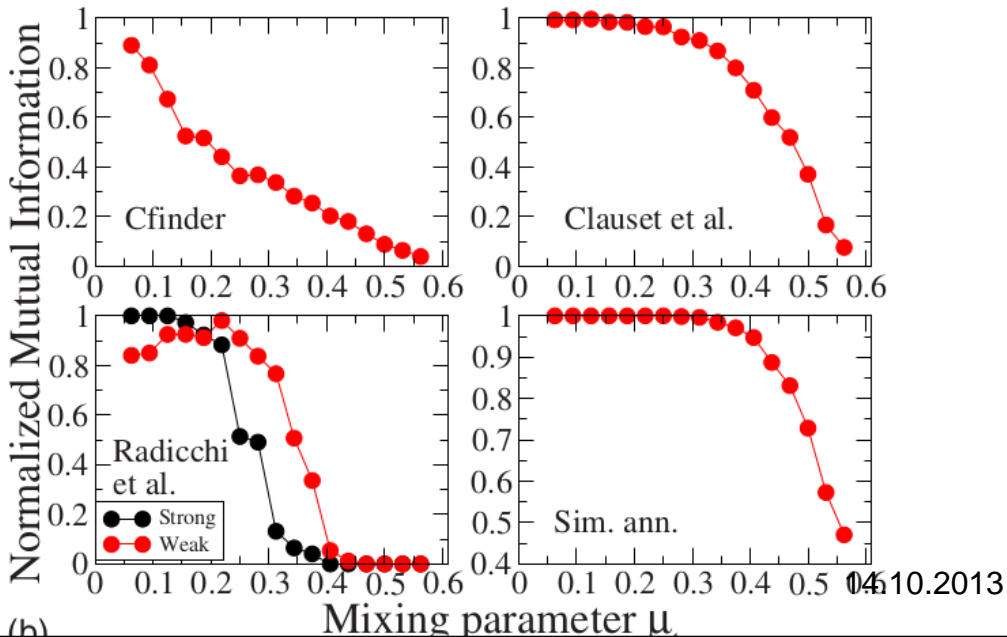
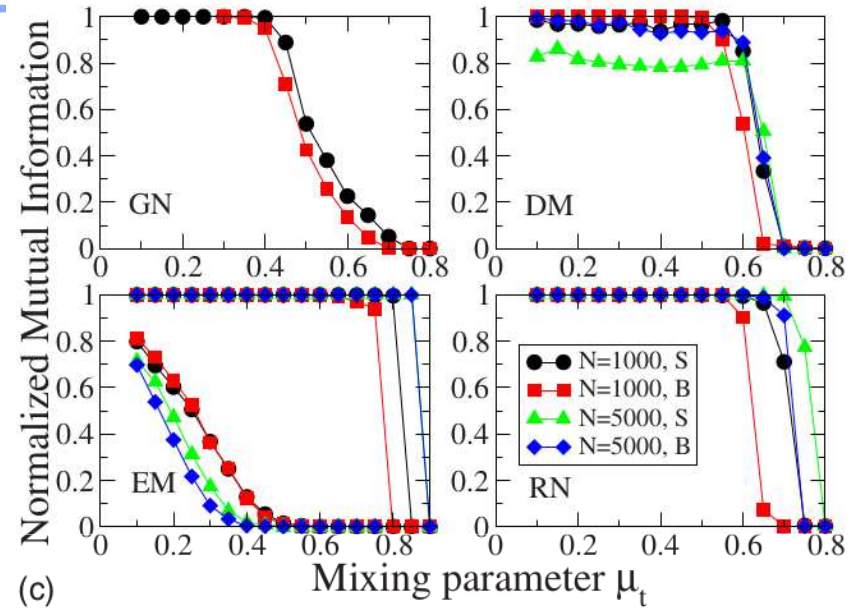
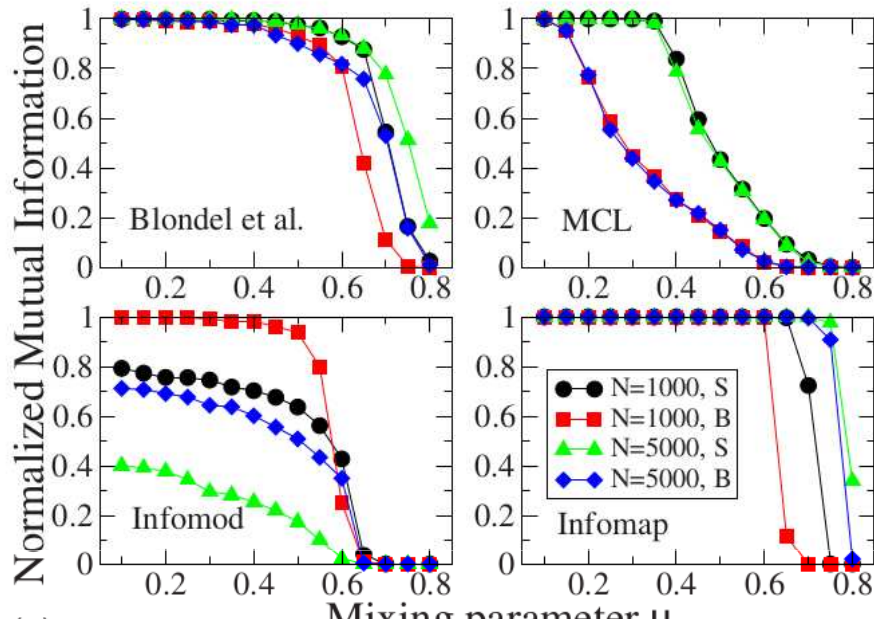


Porovnanie algoritmov na detekciu komunití





Porovnanie algoritmov na detekciu komunití



04.10.2013

(b)

(b)



- Pamäťové obmedzenie
 - Výpočty nad grafmi sú dátovo orientované
 - Prechod topológiou grafu
 - Práca s grafovými dátami je charakteristická náhodným prístupom k dátam
 - Práca s diskom vyžaduje veľa drahých I/O operácií
 - Štandardný prístup: celý graf v pamäti
- Väčšina algoritmov na detekciu má polynomiálnu časovú náročnosť
- Analýza veľkých sietí je problematická
- Príklad:
 - Analýza grafu liniek Wikipédie
 - 3+ miliónov vrcholov, 120+ miliónov hrán
 - Výpočet príliš dlhý
- Pre veľké siete sú výhodné metódy s lineárnou časovou zložitou

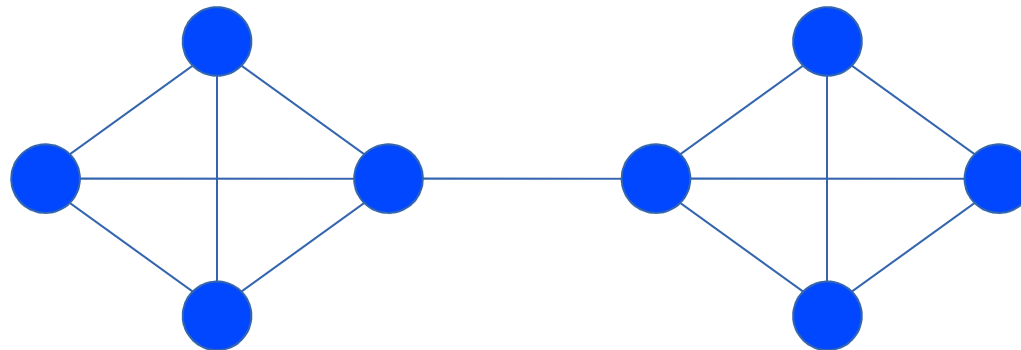


- Greedy algoritmus s pseudo-lineárnou časovou náročnosťou
- Pseudo-lineárna časová náročnosť:
 - Výpočet prebieha v iteráciách
 - Každá iterácia má lineárnu časovú zložitosť
 - Riešenie dobre konverguje po malom počte iterácií => Iterácií stačí relatívne malý počet
- Algoritmus:
 - Inicializácia: každý vrchol v separátnej komunite
 - Iterácie
 - Prechod po uzloch v náhodnom poradí
 - Spracovaný uzol priradíme do komunity, ku ktorej ho viaže najväčší počet hrán
 - V prípade viacerých komunit s najvyšším skóre, vyber príslušnosť náhodne



Európska únia
Európsky fond regionálneho rozvoja

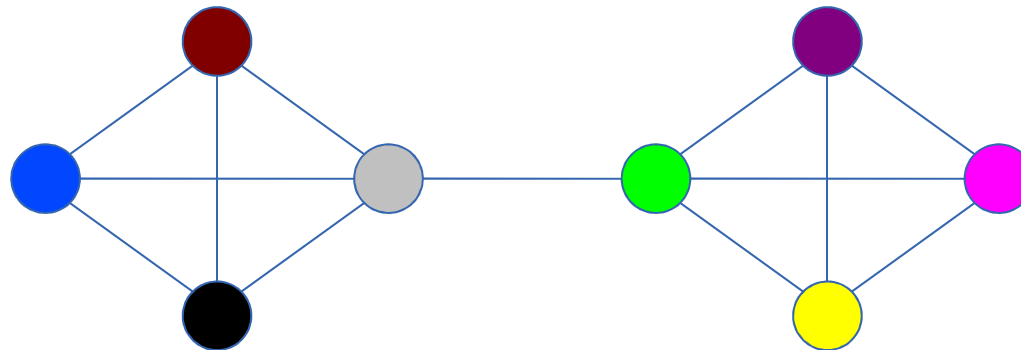
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

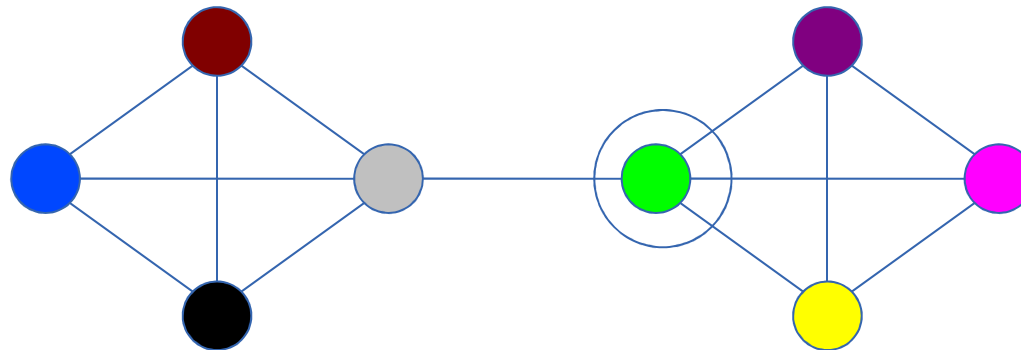
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

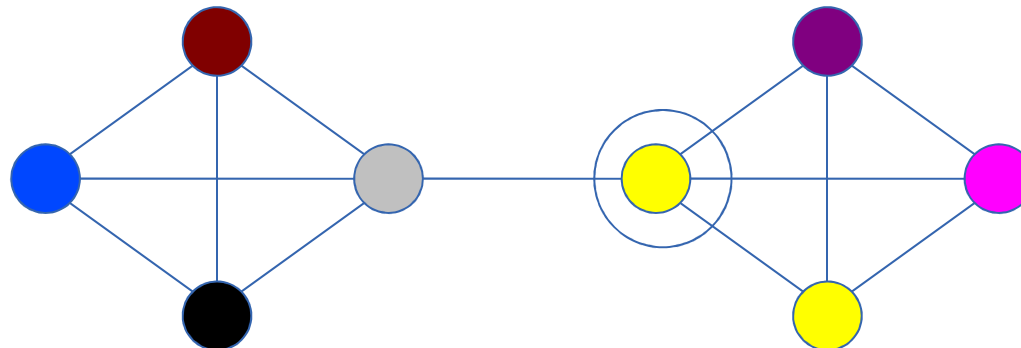
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

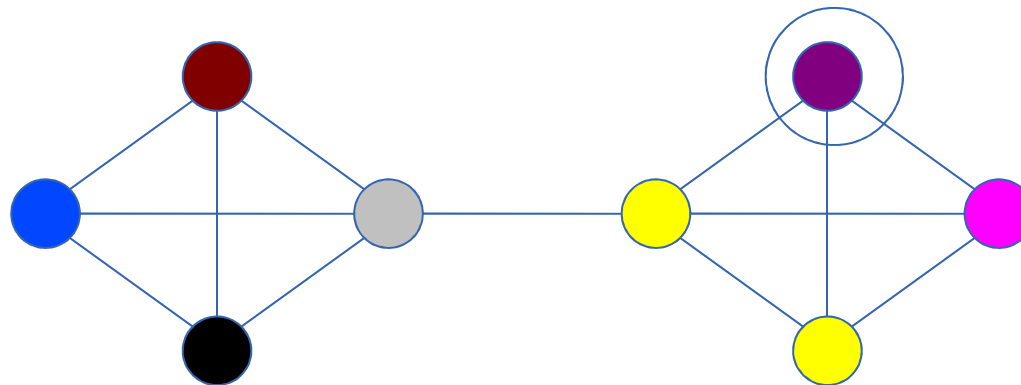
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

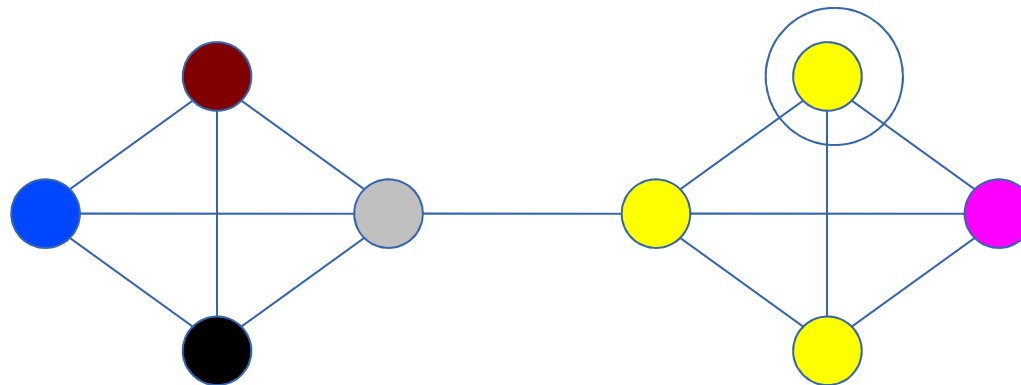
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

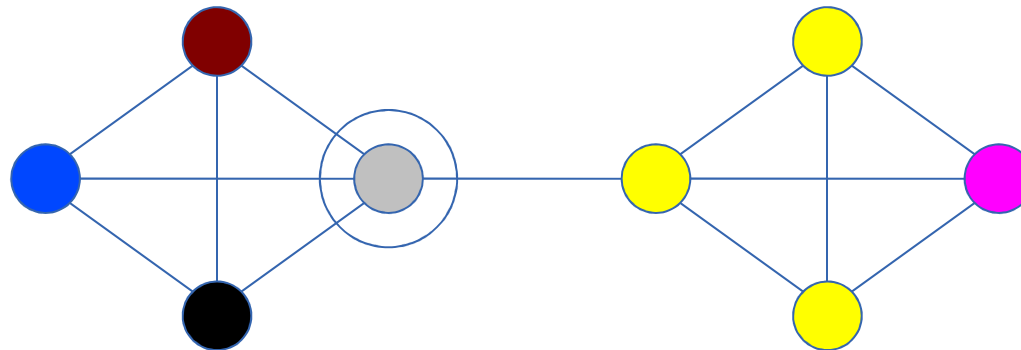
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

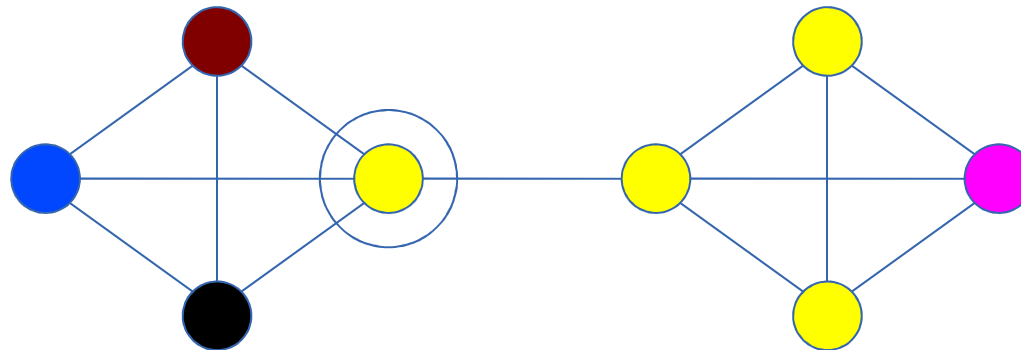
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

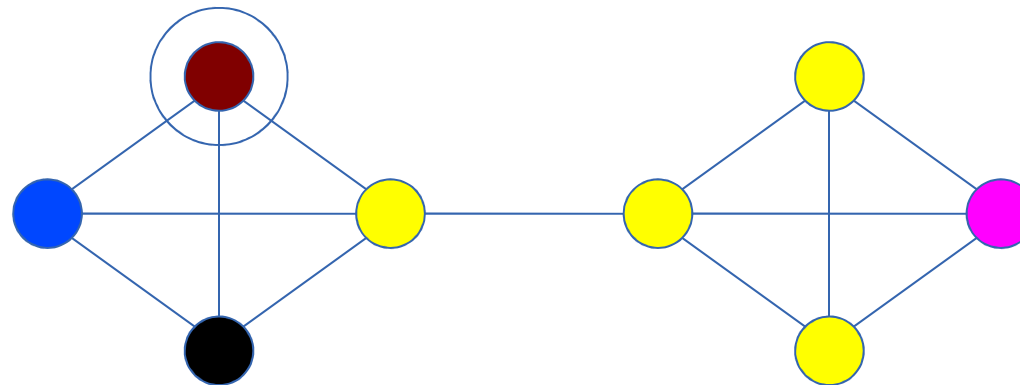
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

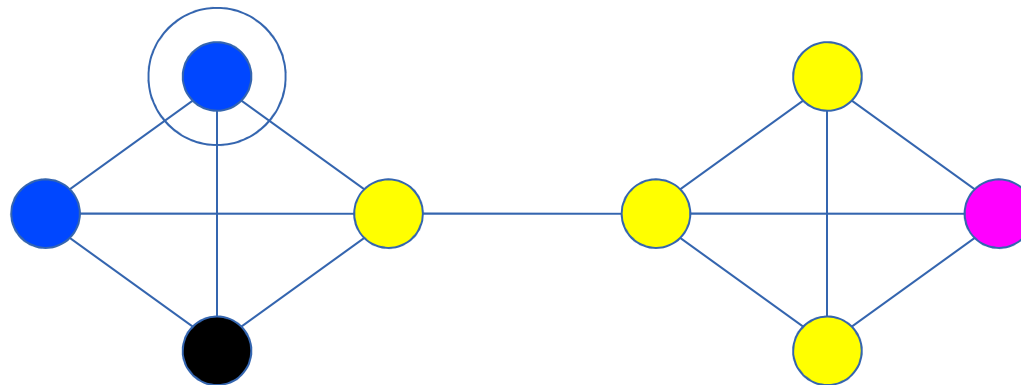
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

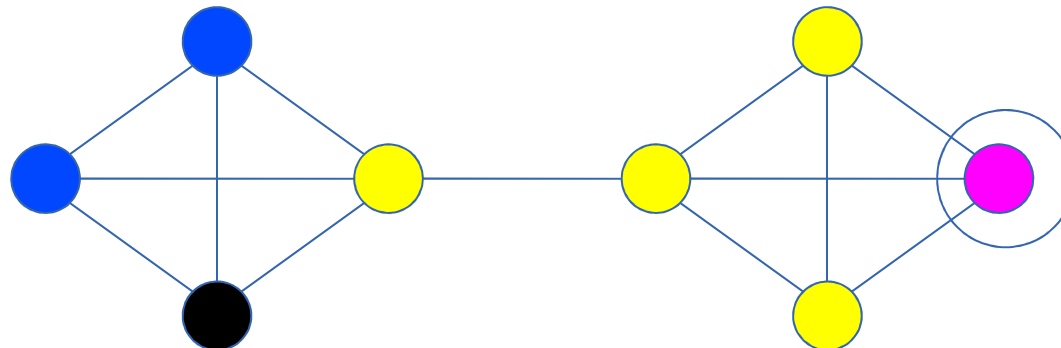
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

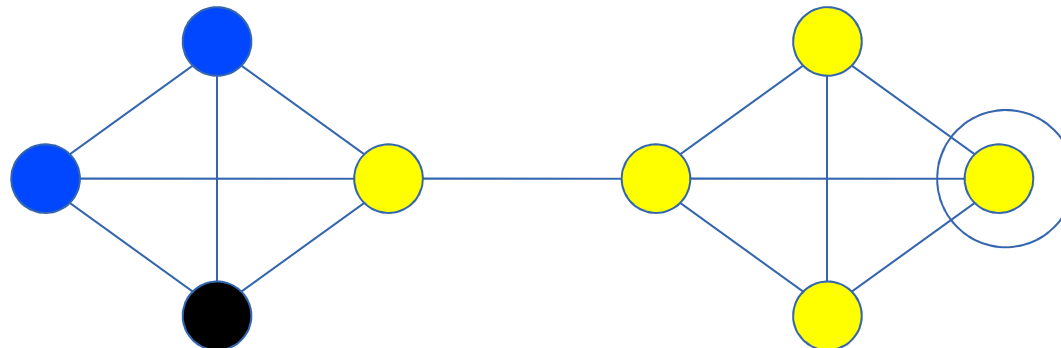
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

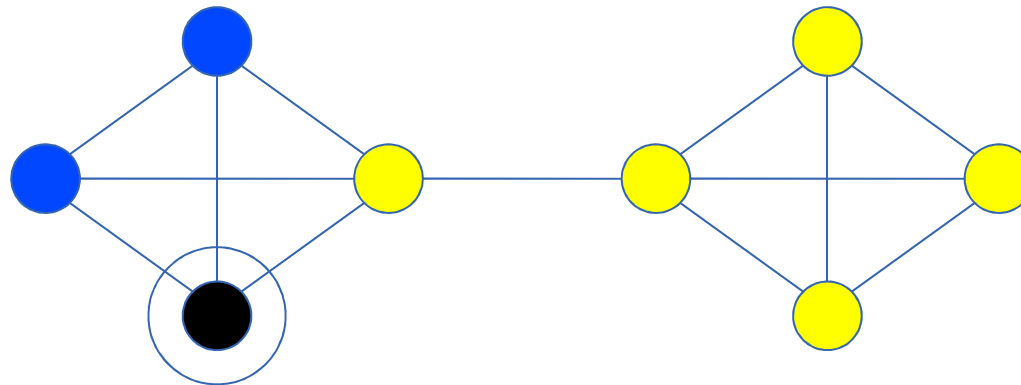
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

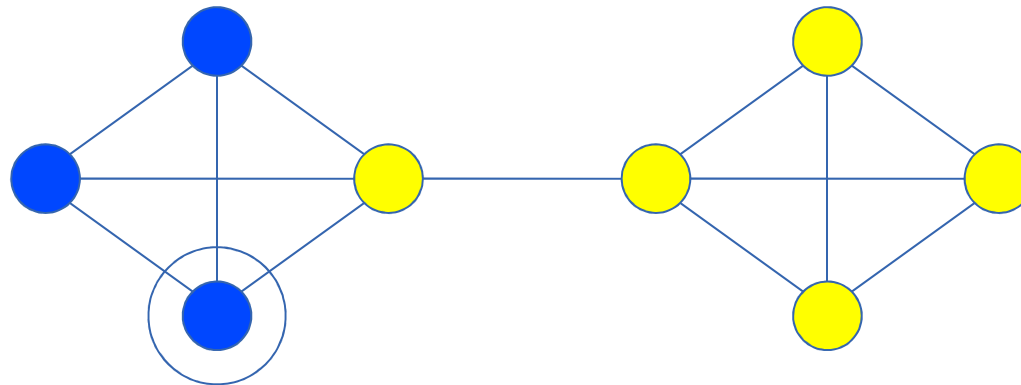
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

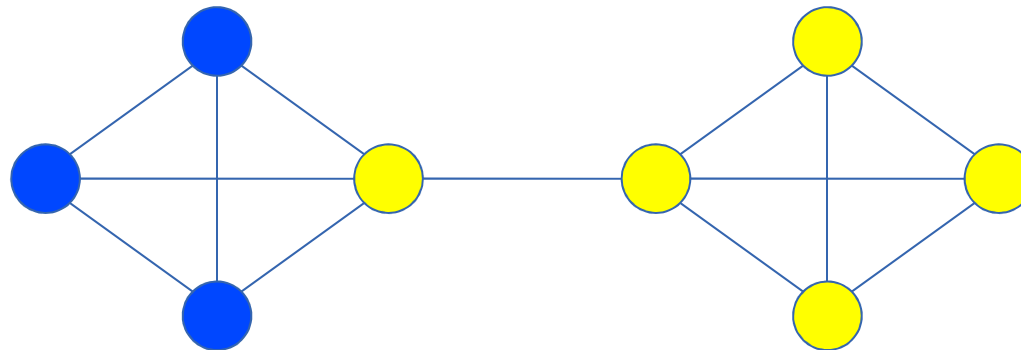
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

Propagácia značiek – Jednoduchý príklad

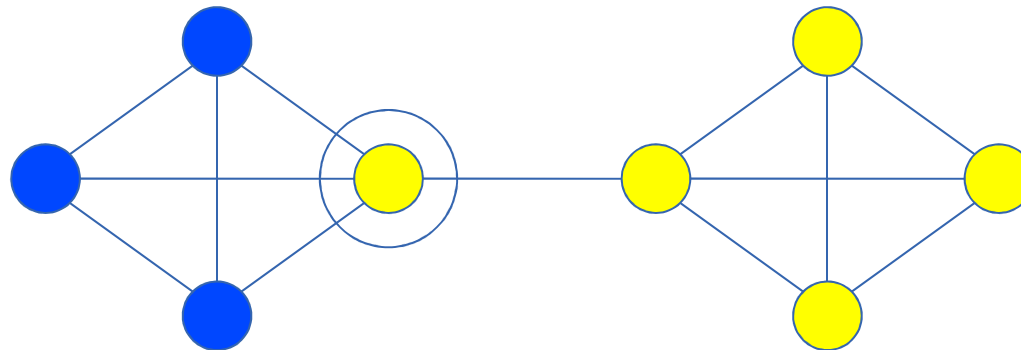


End of 1st iteration



Európska únia
Európsky fond regionálneho rozvoja

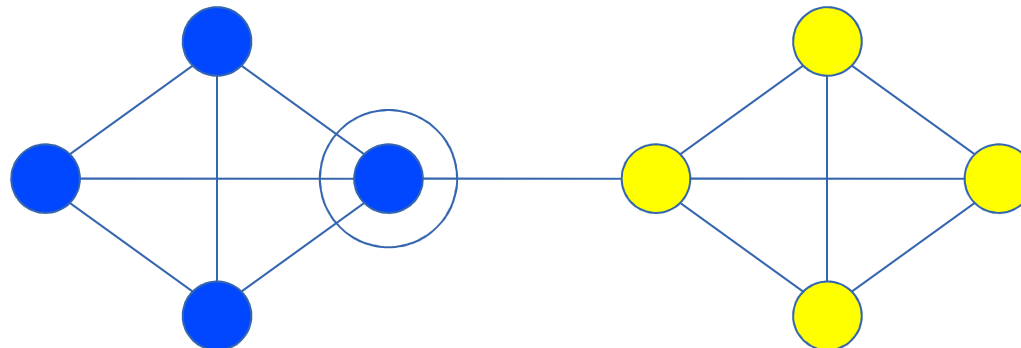
Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

Propagácia značiek – Jednoduchý príklad





Európska únia
Európsky fond regionálneho rozvoja

Propagácia značiek



- Pre benchmark siete s nastavenou vysokou „viditeľnosti komunit“ má dobré výsledky
- Pri menej jasných konfiguráciách má tendenciu zahrnúť veľkú časť vrcholov do jednej komunity



- Základný algoritmus:
 - Inicializácia: každý vrchol v separátnej komunite
 - Iterácie
 - Prechod po uzloch v náhodnom poradí
 - Spracovávaný uzol priradíme do komunity, pre ktorá má najvyššie hodnotu **funkcie nárastu**
 - V prípade viacerých komunit s najvyšším skóre, vyber príslušnosť náhodne



Európska únia
Európsky fond regionálneho rozvoja

Louvain metóda



- Greedy algoritmus optimalizujúci modularitu
- Založený na princípe propagácie značiek
- Zachováva pseudo-lineárnu časovú závislosť
- Funkcia nárastu:

$$\Delta Q = \left[\frac{\sum_{in} + k_{i,in}}{2m} - \left(\frac{\sum_{tot} + k_i}{2m} \right)^2 \right] - \left[\frac{\sum_{in}}{2m} - \left(\frac{\sum_{tot}}{2m} \right)^2 - \left(\frac{k_i}{2m} \right)^2 \right]$$

- SUM_in – suma váh liniek v komunite C
- SUM_tot - suma váh liniek hrán incidentých s uzlami komunity C
- k_i – váha liniek incidentých s uzlom k
- k_i,in – váha liniek incidentých s uzlom k v C
- m - váha všetkých liniek v sieti

[Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre: Fast unfolding of communities in large networks, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008]

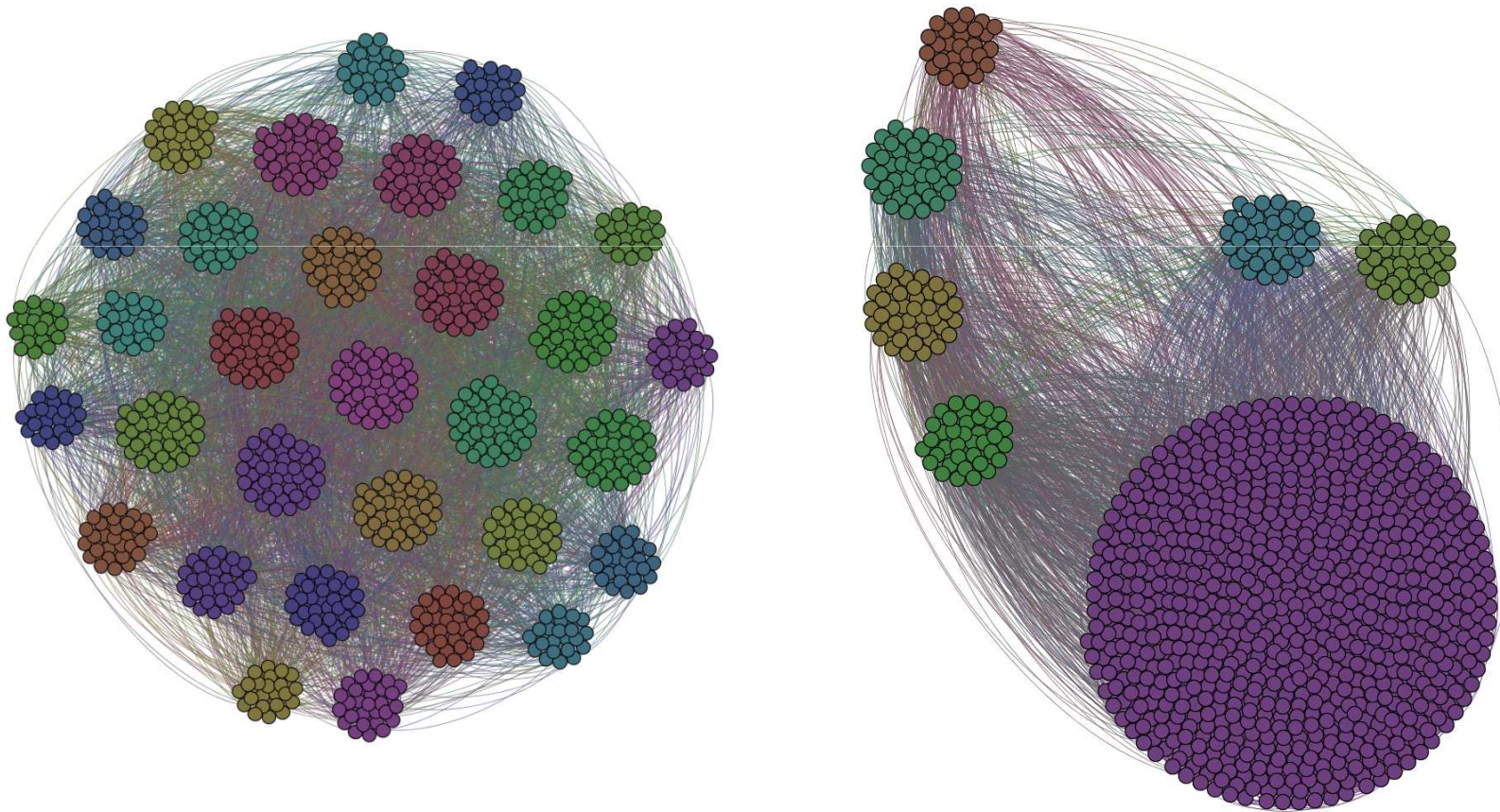


- Detekcia komunít v link grafe Wikipédie
 - 3.1M nodes
 - 91M edges
- Label Propagation algoritmus:
 - Najväčšia kominta: 2.96M uzlov
- Louvain metóda:
 - Najväčšia kominta: 400K uzlov
 - 20 najväčších komunít : 95% siete



Európska únia
Európsky fond regionálneho rozvoja

Cvičenie: Greedy detekcia komúní na Wikipédii





- Greedy algoritmus založený na princípe propagácie značiek
- Zachováva pseudo-lineárnu časovú závislosť
- Motivácia:
 - Zabrániť vzniku obrovských zhlukov
 - Umožniť parametrizáciu maximálnej veľkosti komunity
- Funkcia nárastu:

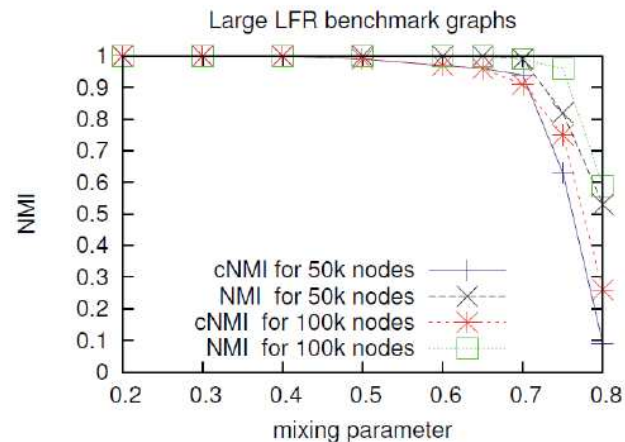
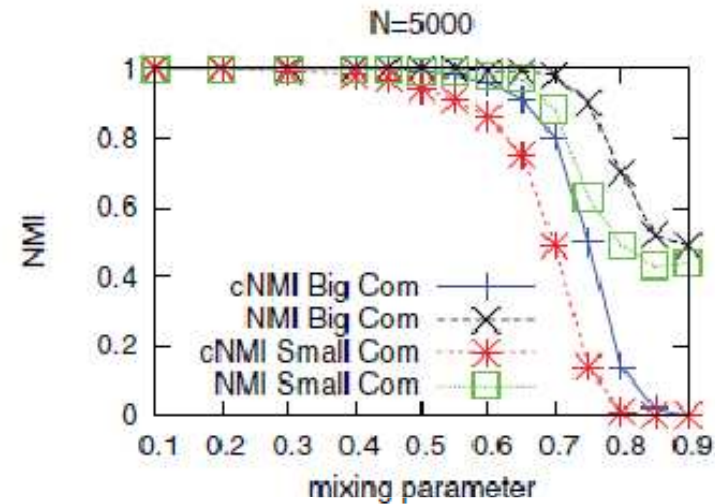
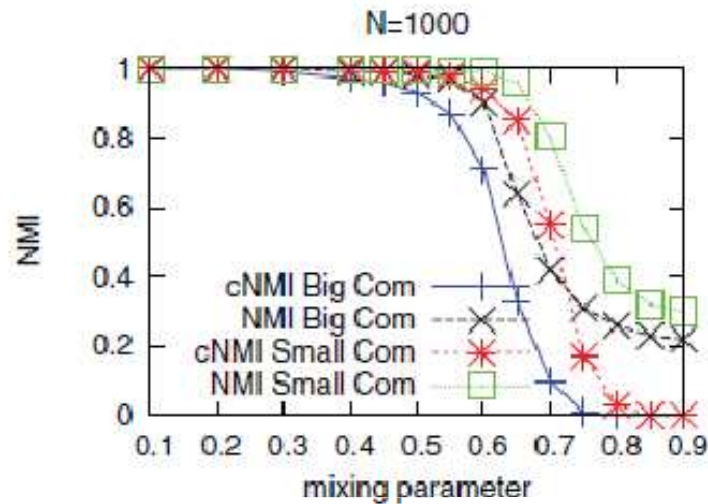
$$seed_gain(n, C) = aff(n, C) \times \log \left(\frac{UpperLimit}{|C|} \right)$$

[Marek Ciglan , Kjetil Nørvgå: Fast detection of size-constrained communities in large networks, proceedings of WISE'10, LNCS Volume 6488/2010, Springer-Verlag]



Európska únia
Európsky fond regionálneho rozvoja

SCCD (Size Constrained Community Detection)



[Marek Ciglan , Kjetil Nørnvåg: Fast detection of size-constrained communities in large networks, proceedings of WISE'10, LNCS Volume 6488/2010, Springer-Verlag]



- Práca: *[J. Yang and J. Leskovec. Defining and Evaluating Network Communities based on Ground-Truth. In: ICDM, 2012]*
- Autori použili reálne siete s explicitne definovanými komunitami
- Siete z rôznych zdrojov:
 - **LiveJournal, Orkut, Youtube**
 - užívateľmi definované skupiny budeme pokladať za komunity
 - **DBLP**
 - Graf ko-autorstva publikácií
 - Autori publikujúci na rovnakých konferenciách a spojený v komponente => komunita
 - **Amazon**
 - Sieť s hranami indikujúcimi časté spolu-kupovanie dvoch produktov
 - Komunity definované kategóriami
- Filtrovanie explicitne definovaných komunit na základe lokálnych metrick kvality komunit



- Práca: *[Marek Ciglan, Michal Laclavík and Kjetil Nørnvåg: On Community Detection in Real-World Networks and the Importance of Degree Assortativity, 19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2013]*
- Použili sme rýchle algoritmy na detekciu komunit na niekoľkých rozsiahlych sieťach reálneho sveta so známymi explicitnými komunitami
 - LiveJournal, Orkut, Youtube, DBLP, Amazon, DBPedia
 - DBPedia: znalostná báza derivovaná z Wikipédie; explicitné komunity sú kategórie Wikipédie
 - Použité algoritmy: Label propagation, Louvain method, SCCD
- Pre niektoré siete algoritmy doručili dobrú aproximáciu explicitne definovaných zhlukov, pre časť sietí nie



- Porovnanie detekovaného rozdelenia siete a explicitných komunit

	Louvain	Label Prop.	SCCD
DBLP-top5K	0.53 (0.24)	0.49 (0.25)	0.51 (0.26)
Amazon-top5K	0.76 (0.37)	0.85 (0.46)	0.86 (0.46)
YouTube-top5K	0.23 (0.09)	0.08 (0.02)	0.19 (0.05)
Orkut-top5K	0.04 (0.03)	0.03 (0.02)	0.24 (0.06)
LJ-top5K	0.52 (0.28)	0.57 (0.32)	0.60 (0.32)
DBPedia-top5K	0.004 (0.0008)	0.003 (0.003)	0.13 (0.06)



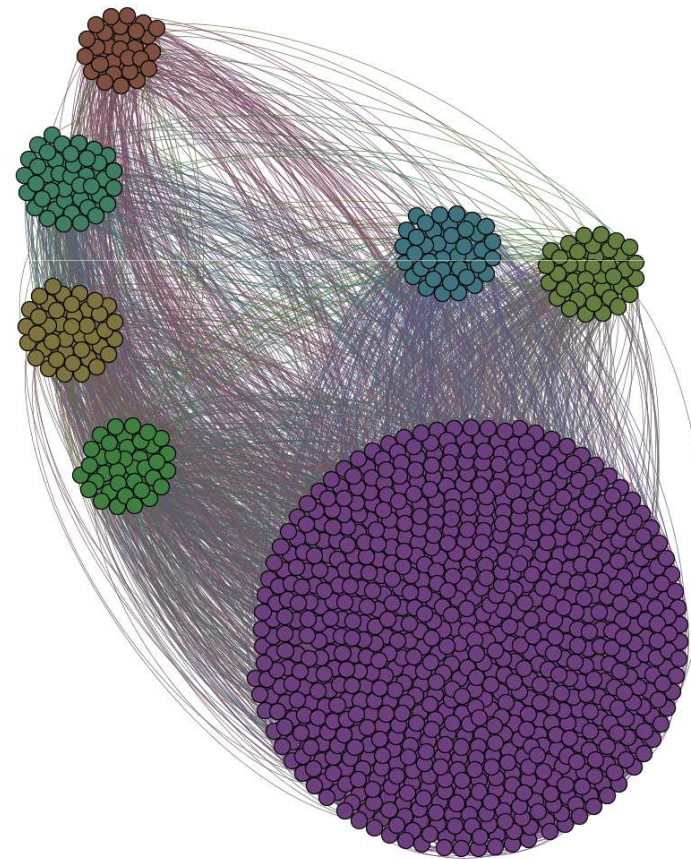
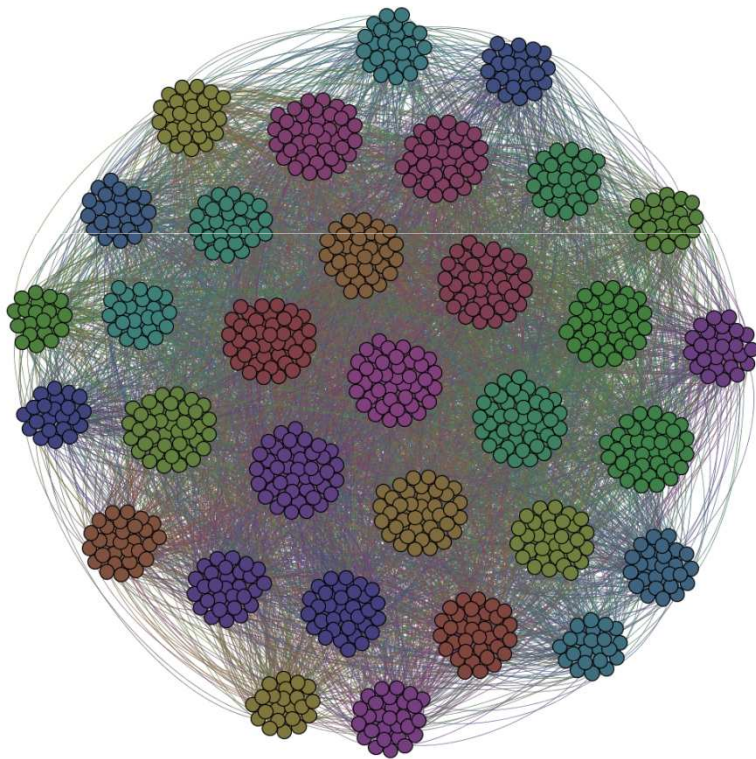
- Časť vrcholov prislúchajúcich 50-tim najväčším komunitám

	Lou fract. of N	LP fract. of N	SCCD fract. of N	Ground-truth fract. of N
DBLP	0.11	0.07	0.05	0.32
Amazon	0.03	0.03	0.04	0.83
YouTube	0.50	0.66	0.17	0.01
LJ	0.67	0.71	0.07	0.09
Orkut	0.97	0.99	0.93	0.09
DBPedia	0.98	0.96	0.89	0.001



Európska únia
Európsky fond regionálneho rozvoja

Detekcia komunití v rozsiahlych sieťach reálneho sveta





- **Hypotéza:** tendencia skúmaných algoritmov založených na propagácii značiek na sieťach DBPedia a Orkut zhrnúť väčšinu vrcholov do malého počtu veľkých komunit je spôsobená veľkým podielom hrán spájajúcich väčšinu vrcholov s vrcholmi vysokého stupňa v hustom jadre siete.
- **Koeficient asortativity siete:** číselne vyjadruje tendenciu uzlov byť spojených s uzlami podobného stupňa (Pearsonov korelačný koeficient stupňov vrcholov spojených hranou)



- Asortativita skúmaných sietí a assortativita ich štruktúry explicitných komunití

	net AC	Top5k Comm. AC	All Comm. AC
DBLP	0.267	0.436	0.446
Amazon	-0.059	-0.077	-0.026
YouTube	-0.037	0.067	0.068
LJ	0.045	0.464	0.365
Orkut	0.016	0.233	0.326
DBPedia	-0.018	0.958	0.973

- Zistenie: Asortativita sietí a štruktúry komunití môže byť veľmi rozdielna.



- Skúmali sme možnosť modifikovať sieť pridaním váh na hrany
- Cieľom bolo znížiť váhu liniek spájajúcich disasortatívne uzly

Weighting function 1:

$$f_1(z, y) = \text{floor}(\log_{10}(\max(z, y)/\min(z, y)))$$

$$w1(e_{i,j}) = 10^{-f_1(d(i), d(j))}$$

Weighting function 2:

$$f_2(z, y) = \text{floor}(\log_{10}(|z - y|))$$

$$w2(e_{i,j}) = 10^{-f_2(d(i), d(j))}$$



Network	Alg.	Orig	w1	w2
DBLP	Lou	0.54 (0.24)	0.54 (0.26)	0.56 (0.26)
	LP	0.49 (0.25)	0.49 (0.24)	0.53 (0.25)
	SCCD	0.52 (0.26)	0.51 (0.25)	0.54 (0.25)
Amazon	Lou	0.76 (0.37)	0.84 (0.43)	0.83 (0.42)
	LP	0.85 (0.46)	0.85 (0.46)	0.83 (0.44)
	SCCD	0.86 (0.46)	0.86 (0.46)	0.83 (0.44)
YouTube	Lou	0.23 (0.10)	0.23 (0.11)	0.31 (0.12)
	LP	0.08 (0.02)	0.14 (0.05)	0.24 (0.09)
	SCCD	0.19 (0.06)	0.21 (0.08)	0.26 (0.10)
LJ	Lou	0.52 (0.28)	0.49(0.27)	0.52 (0.28)
	LP	0.57 (0.32)	0.59 (0.32)	0.62 (0.32)
	SCCD	0.6 (0.32)	0.61 (0.32)	0.62 (0.33)
Orkut	Lou	0.04 (0.03)	0.04 (0.04)	0.17 (0.05)
	LP	0.03 (0.02)	0.03 (0.03)	0.11 (0.04)
	SCCD	0.24 (0.06)	0.20 (0.06)	0.25 (0.06)
DBPedia	Lou	0.004 (0.0008)	0.016 (0.03)	0.053 (0.15)
	LP	0.003 (0.003)	0.054 (0.14)	0.31(0.17)
	SCCD	0.13 (0.06)	0.26 (0.15)	0.34 (0.18)



Európska únia
Európsky fond regionálneho rozvoja

Detekcia komunití v rozsiahlych sieťach reálneho sveta

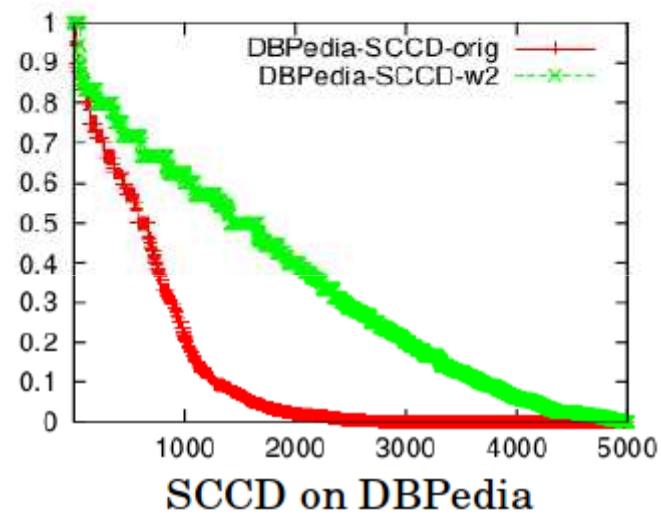
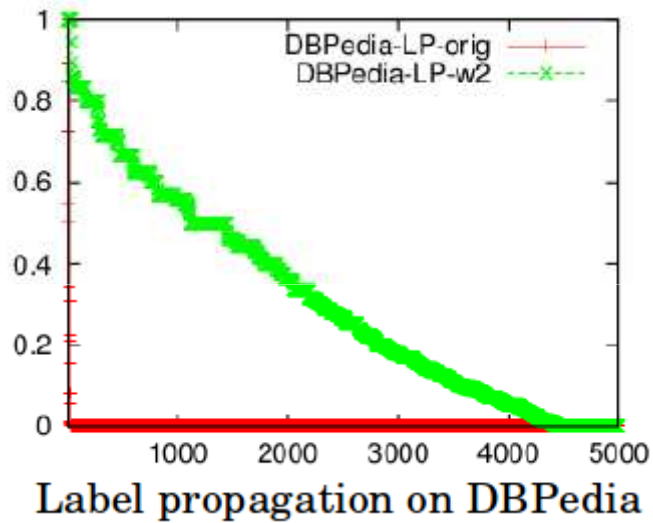


Net	Alg.	w_2	rnd	MW	$MW \times w_2$
DBLP	Lou	0.56	0.47(0.23)	0.54 (0.26)	0.56(0.26)
	LP	0.53	0.48(0.23)	0.56 (0.27)	0.56(0.26)
	SCCD	0.54	0.52(0.24)	0.57 (0.27)	0.57(0.27)
Amazon	Lou	0.83	0.83(0.43)	0.85 (0.43)	0.83(0.42)
	LP	0.83	0.82(0.43)	0.85 (0.44)	0.84(0.43)
	SCCD	0.83	0.82(0.42)	0.85 (0.44)	0.83(0.42)
YouTube	Lou	0.31	0.15(0.05)	0.20 (0.09)	0.30(0.12)
	LP	0.24	0.11(0.04)	0.17 (0.07)	0.29(0.11)
	SCCD	0.26	0.21(0.06)	0.32 (0.13)	0.32(0.12)
LJ	Lou	0.52	0.50(0.27)	0.53 (0.28)	0.53(0.29)
	LP	0.62	0.56(0.31)	0.61 (0.33)	0.63(0.33)
	SCCD	0.62	0.62(0.33)	0.63 (0.34)	0.63(0.33)
Orkut	Lou	0.17	0.06(0.03)	0.08(0.03)	0.13(0.04)
	LP	0.11	0.03(0.02)	0.06(0.03)	0.18(0.04)
	SCCD	0.25	0.22(0.05)	0.16 (0.04)	0.19(0.05)
DBPedia	Lou	0.05	0.03(0.01)	0.05(0.03)	0.29(0.15)
	LP	0.31	0.001 (0.001)	0.33 (0.16)	0.40 (0.20)
	SCCD	0.34	0.15(0.06)	0.36 (0.18)	0.41(0.20)



Európska únia
Európsky fond regionálneho rozvoja

Detekcia komunití v rozsiahlych sieťach reálneho sveta





Európska únia
Európsky fond regionálneho rozvoja

Výskumné smery



- Rozširovanie základného výskumu v algoritmoch
- Aplikácie detekcie komunit
- Škálovateľnosť



Európska únia
Európsky fond regionálneho rozvoja



Aplikácie detekcie komunít



Európska únia
Európsky fond regionálneho rozvoja

Aplikácie



- Analýza sociálnych sietí
 - Rozdelenie do skupín
- WWW – analýza liniek
 - skupiny tématicky súvisiacich stránok
- Biologické siete
 - Interakčné siete proteínov – komunity zoskupujú proteíny rovnakej funkcie
 - Funkčné moduly v metabolických
- Business intelligence
 - Teleco – Churn prediction
 - Finančný sektor – fraud detection

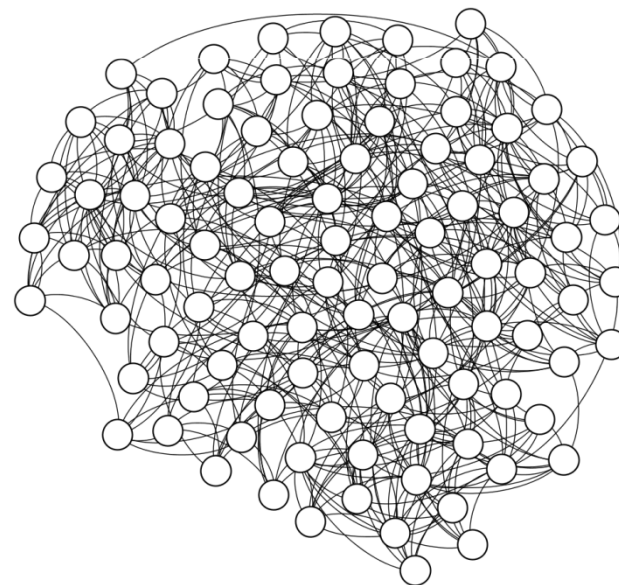


- Churn prediction
 - Cieľ je predikovať možný odchod zákazníkov od mobilného operátora
 - Udržanie zákazníka lacnejšie ako získanie nového

- Prediktívna analýza
 - Dáta o zákazníkoch – vektory rysov
 - $(f_1, f_2, \dots, f_n \mid \text{churn}=?)$



- Churn prediction
 - Cieľ je predikovať možný odchod zákazníkov od mobilného operátora
 - Udržanie zákazníka lacnejšie ako získanie nového
- Prediktívna analýza
 - Dáta o zákazníkoch – vektory rysov
 - $(f_1, f_2, \dots, f_n \mid \text{churn}=?)$
 - Rozšírenie vektorov o grafové rysy
 - Príslušnosť do komunity
 - Centralita uzla (napr. page rank)
 -
 - $(f_1, f_2, \dots, f_n, gf_1, gf_2, \dots, gf_m \mid \text{churn}=?)$





- **Motivácia**

- **Sociálna sieť – relácie medzi užívateľmi**

- Facebook - 750 miliónov užívateľov
 - YouTube - 490 miliónov pravidelných užívateľov
 - Twitter - 550 miliónov užívateľov
 - Wikipedia - 91,000 kontribútorov

- **Obsah generovaný užívateľmi**

- textové správy, fotky, videá, reakcie (+1 / Likes)
 - Facebook - 30 miliárd zdieľaných položiek / mesiac
 - Twitter - 190 miliónov mikropostov / deň
 - Wikipedia - 17 miliónov hostovaných článkov

- **Interakcia užívateľov s obsahom**

- YouTube - 92 miliárd zobrazení stránok / mesiac
 - Twitter - 1.6 miliárd dopytov za deň
 - Facebook – priemerný čas strávený na stránke za mesiac: 15H 33M



- Výzvy:
 - Časová náročnosť algoritmov
 - Pamäťové ohraničenie
- Pamäťové ohraničenie
 - Distribuované rámce pre výpočty nad grafmy
 - Koncept „Pregel“ od googlu

[Grzegorz Malewicz, Matthew H. Austern, Aart J.C Bik, James C. Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. 2010. Pregel: a system for large-scale graph processing. In Proceedings of the 2010 ACM SIGMOD International Conference on Management of data (SIGMOD '10). ACM, New York, NY, USA, 135-146. DOI=10.1145/1807167.1807184]

- Open source implementácie
 - Giraph, Apache Hama GraphLab



Európska únia
Európsky fond regionálneho rozvoja



Applications of Graphs and Networks



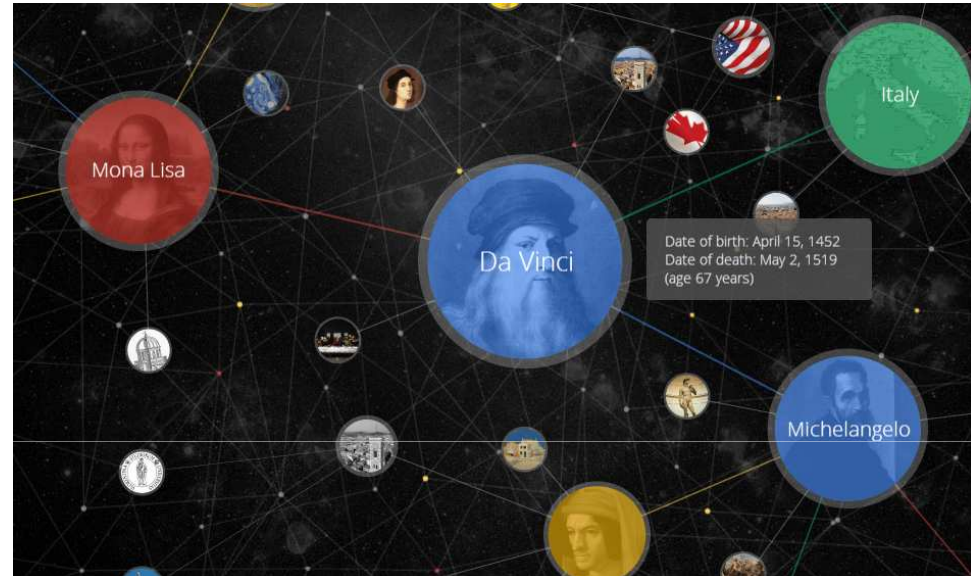
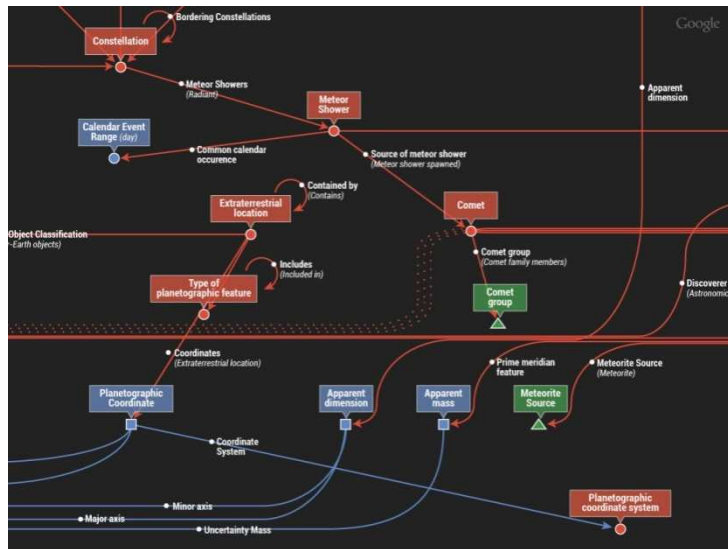
Európska únia
Európsky fond regionálneho rozvoja

Google Knowledge Graph



[ulanoff]

- Wikipedia
- Freebase
- Confirmed human knowledge



Mona Lisa



The Mona Lisa is a half-length portrait of a woman by the Italian artist Leonardo da Vinci, which has been acclaimed as "the best known, the most visited, the most written about, the most sung about, the most parodied work of art in the world." [Wikipedia](#)

Started: 1503
Completed: 1505
Location: Louvre
Dimensions: 30.3" x 20.9" (77 cm x 53 cm)
Genre: Portrait
Media: Oil paint

Leonardo da Vinci



Leonardo di ser Piero da Vinci was an Italian Renaissance polymath: painter, sculptor, architect, musician, scientist, mathematician, engineer, inventor, anatomist, geologist, cartographer, botanist, ... [Read more on en.wikipedia.org](#)

Born: April 15, 1453, Anchiano
Died: May 2, 1519, Clos Lucé
Buried: St. Florentin's Church
Inventions: Viola organista, Double hull
Parents: Caterina da Vinci, Piero da Vinci

Works



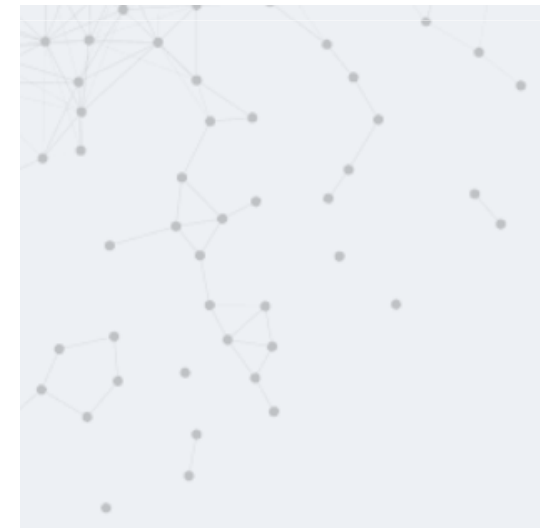
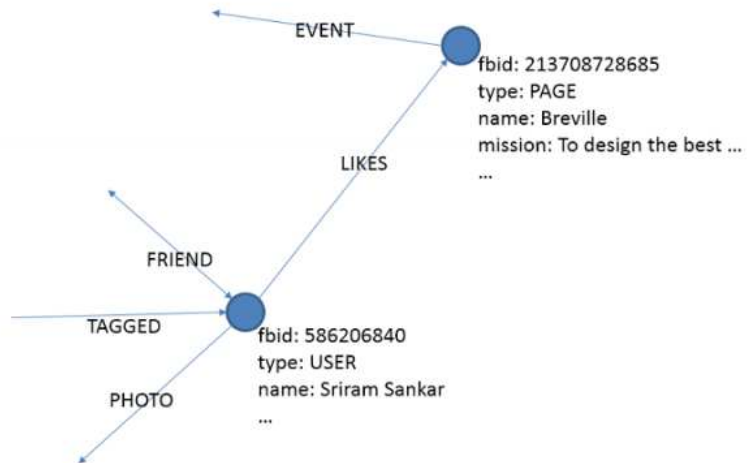
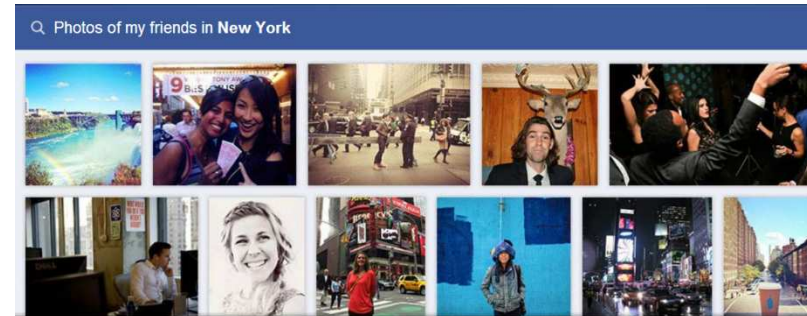
People also search for





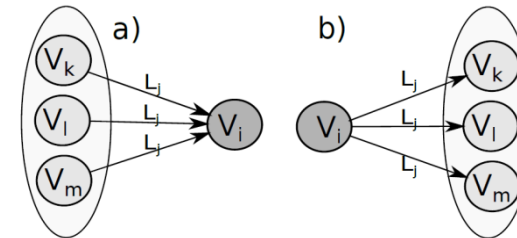
[facebook13]

- Užívateľmi generovaný obsah
- Prepojenia na web

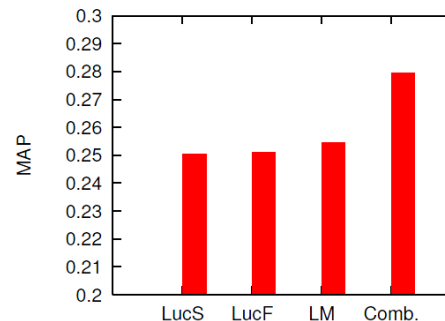
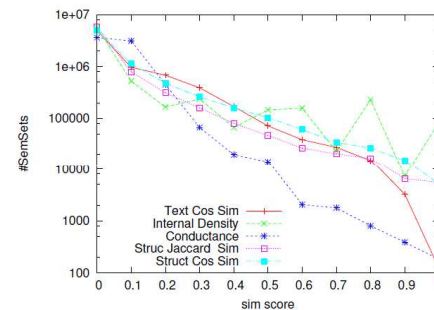




- Answering list type questions: *astronauts who walked on the Moon*
- Wikipedia as text and networks/graph
- Text: IR methods, Lucene based
- Graph/network: spreading activation and SemSets
- Winning solution on Semantic Search Challenge 2011



[SemSets]



1. Eugene_Cernan
2. Alan_Bean
3. David_Scott
4. John_Young_(astronaut)
5. Neil_Armstrong
6. Pete_Conrad
7. Harrison_Schmitt
8. Alan_Shepard
9. Charles_Duke
10. Buzz_Aldrin
11. James_Irwin
12. Edgar_Mitchell



- Motivácia

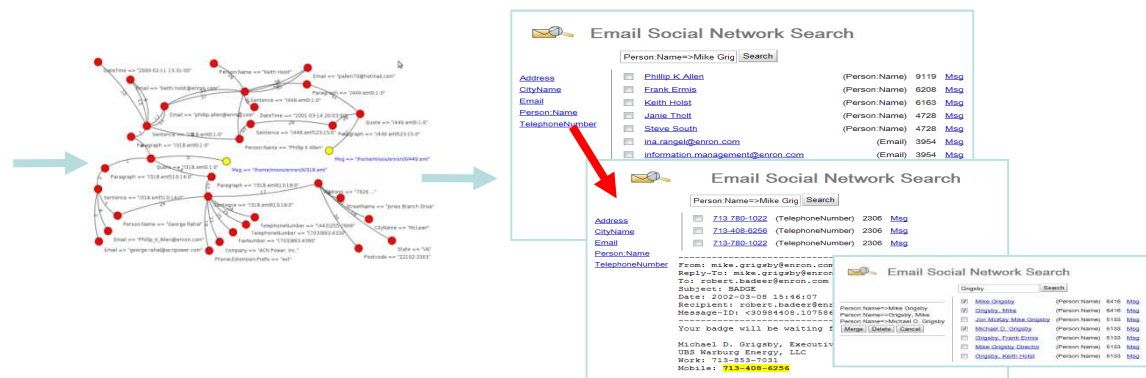
- Grafy a siete sú všadeprítomné : sociálne siete, web, LinkedData, transakcie, komunikácia (email, telefóny).
- Text tiež môže byť prevedený na graf.
- Prepojenie grafových dát a vyhľadávania relácií v nich je dôležité

- Prístup

- Tvorba sémantických stromov a grafov z textu, webu, komunikácie, databáz a LinkedData
- Užívateľská interakcia s týmito dátami aby sa dali lepšie integrovať zdroje a vyčistiť/upraviť dáta
- Užívatelia to budú robiť ak to bude mať zmysel, teda okamžitý vplyv na lepšie výsledky vyhľadávania



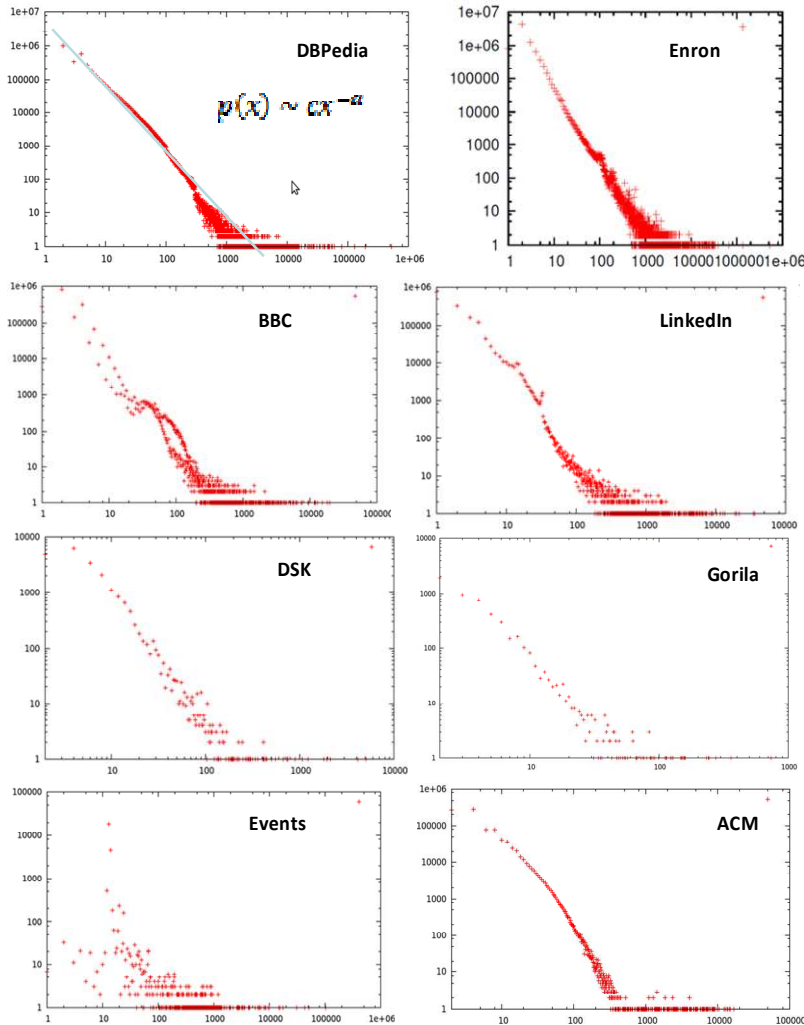
LinkedIn





Európska únia
Európsky fond regionálneho rozvoja

Vlastností vybraných grafov/sietí



Názov siete	Počet vrcholov	Počet hrán	Priem. klást. koef.	Koef. assort.	Priem. najkr. cesta
Enron Full	8 269 278	20 383 709	0,29	-0,02	6,58
Enron5	160 387	630 330	0,30	-0,04	6,64
LinkedIn	1 564 698	6 094 634	0,36	0,13	6,48
BBC	1 725 900	6 839 358	0,34	-0,05	7,55
DSK	21 518	98 952	0,31	0,39	5,79
DSK3	2 857	8 754	0,36	-0,14	5,46
Gorila	5 959	23 724	0,31	0,03	6,25
Events	25 478	539 328	0,38	-0,25	2,47
ACM	941 322	2 198 001	0,34	-0,06	7,30

Datsety:

- DBPedia
- Web
 - BBC, LinkedIn, DSK
- Gorila – document
- Events – agent simulation event graph
- ACM – publications, LinkedData



- ❖ Regex patterns
- ❖ Gazetteers
- ❖ Results
 - ❖ Key-value pairs
 - ❖ Structured into trees
 - ❖ graphs
- ❖ Transformers, Configuration
- ❖ Automatic loading of extractors
- ❖ Visual Annotation Tool
- ❖ Integration with external tools
 - ❖ GATE, Stemers, Hadoop ...
- ❖ Multilingual tests
 - English, Slovak, Spanish, Italian

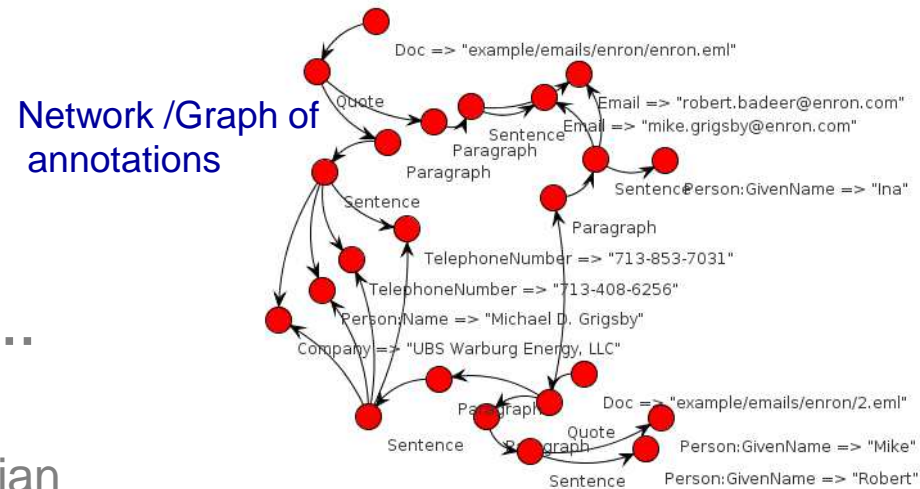
-----ENVELOPE-START-----
 From: mike.grigsby@enron.com
 To: robert.badeer@enron.com
 Subject: BADGE
 Date: 2002-03-08 15:46:07
 -----ENVELOPE-END-----

Your badge will be waiting for you at the front desk in the north tower on mon. if not, then call and we will retrieve you.

Michael D. Grigsby, Executive Director
 UBS Warburg Energy, LLC
 Work: 713-853-7031
 Mobile: 713-408-6256

Text with annotations

Tree of annotations





Európska únia
Európsky fond regionálneho rozvoja

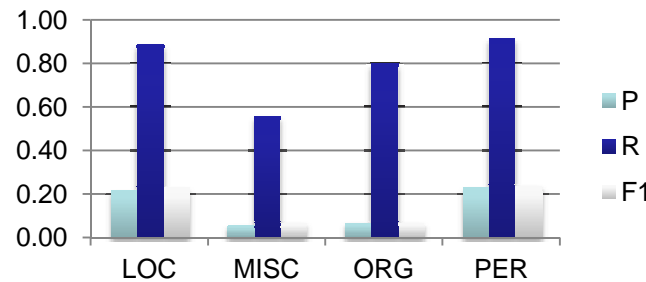
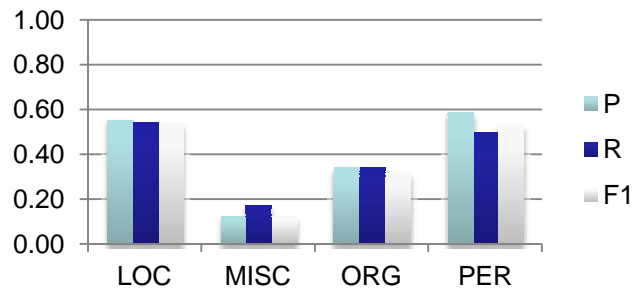
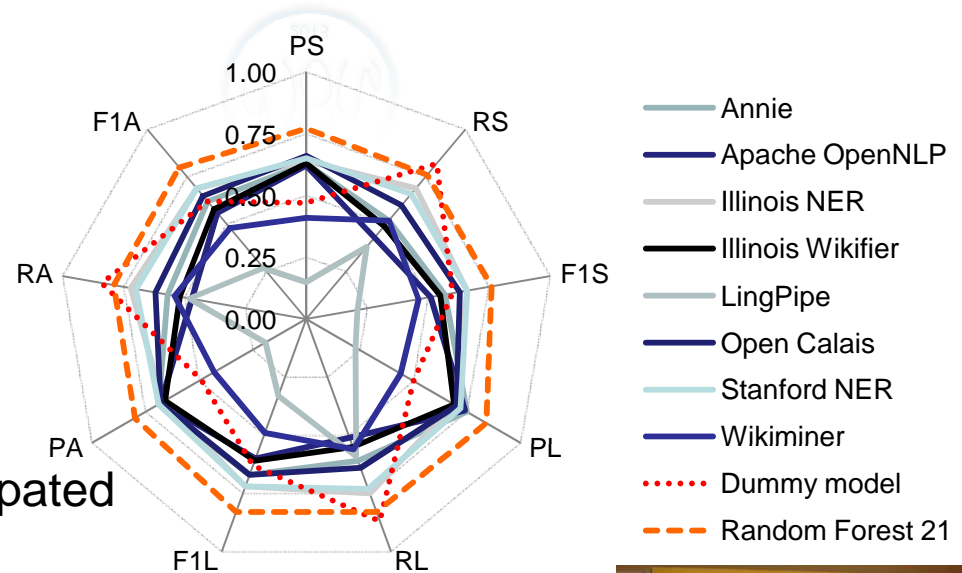
Named Entity Recognition (NER)



- Combination of Existing NER
 - ANNIE (GATE), Apache OpenNLP,
 - Illinois NER, Illinois Wikifier,
 - LingPipe, Open Calais
 - Stanford NER, WikiMiner,
 - Miscinator
- Machine Learning
 - Decision Trees models
- Received second place at MSM 2013, missing first place by 1%, where participated 17 teams world wide

<http://ikt.ui.sav.sk/index.php?n=Main.IEChallenge2013>

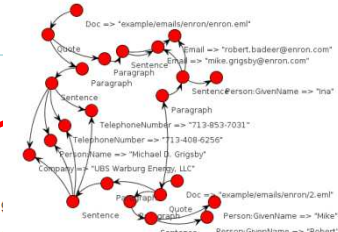
Micro Summary (test set)





- Entity relation search in semantic networks/graphs
- Search, Navigation, Data Interaction
- Aiming at data integration of
 - Structured data (Relational data, LinkedData)
 - Unstructured Data (text, documents, communication)
- Applications:
 - Email, Web, Text documents, LinkedData

<http://ikt.ui.sav.sk/esns/>



Email Social Network Search

Person:Name=>Mike Grig Search

Address
CityName
Email
Person:Name
TelephoneNumber

- Phillip K Allen (Person:Name) 6163 Msg
- Frank Ermis (Person:Name) 4728 Msg
- Keith Holst (Person:Name) 4728 Msg
- Janie Tholt (Email) 3954 Msg
- Steve South (Email) 2054 Msg
- ina_rangel@enron.com
- information.management@enron.com

Email Social Network Search

Grigsby Search

- + 713-408-6256 (TelephoneNumber) 1024 Msg
- + 713-853-7031 (TelephoneNumber) 1000 Msg
- + 713-780-1022 (TelephoneNumber) 24 Msg
- + 713 780-1022 (TelephoneNumber) 24 Msg

TelephoneNumber=>713-408-6256
Neighbor Count:8

Mike Grigsby(Person:Name)
Michael D. Grigsby(Person:Name)
UBS Warburg Energy, LLC(Company)

From: mike.grigsby@enron.com
To: robert.badeer@enron.com
Subject: BADGE
Date: 2002-03-08 15:46:07

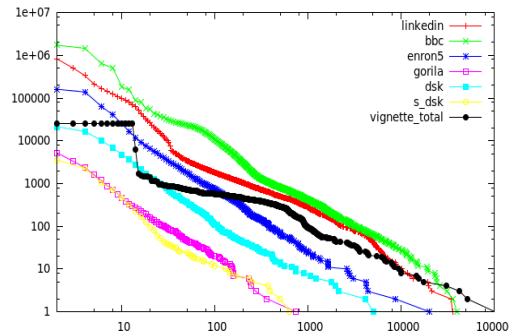
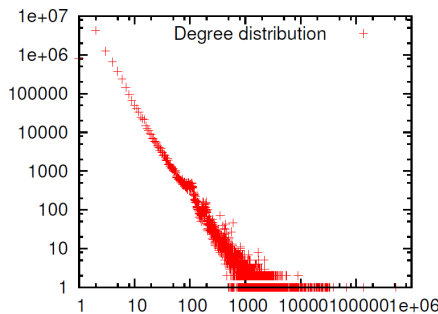
Your badge will be waiting for you at the front desk in the north tower

Michael D. Grigsby, Executive
UBS Warburg Energy, LLC
Work: 713-853-7031
Mobile: 713-408-6256

Email Social Network Search

Grigsby Search

- Mike Grigsby (Person:Name) 6416 Msg
- Grigsby, Mike (Person:Name) 6416 Msg
- Jon McKay, Mike Grigsby (Person:Name) 5133 Msg
- Michael D. Grigsby (Person:Name) 5133 Msg
- Grigsby, Frank Ermis (Person:Name) 5133 Msg
- Mike Grigsby Director (Person:Name) 5133 Msg
- Grigsby, Keith Holst (Person:Name) 5133 Msg

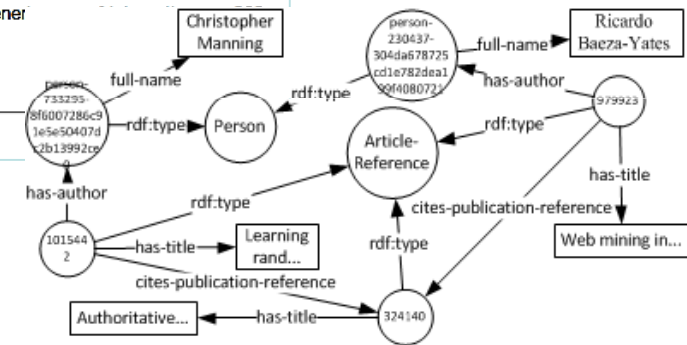




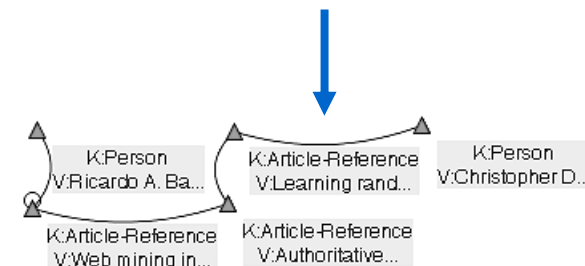
Graph based Semantic Search

<p><u>Paper</u></p> <p><u>addresses-generic-area-of-interest</u></p> <p>Ricardo A. Baeza-Yates(Author) Christopher D. Manning(Author) Ricardo Baeza-Yates(Author) R. Baeza-Yates(Author) R. A. Baeza-Yates(Author) Christopher Manning(Author)</p> <p> <input checked="" type="radio"/> AND <input type="radio"/> OR </p> <input type="button" value="Search Multi"/>	<table border="0" style="width: 100%;"> <tr> <td style="width: 5%; text-align: right;">+</td> <td style="width: 45%;">H.3.3</td> <td style="width: 45%;">(addresses-generic-area-of-interest) 2804692</td> </tr> <tr> <td style="text-align: right;">+</td> <td>The anatomy of a large-scale hypertextual Web search engine</td> <td>(Paper) 122898</td> </tr> <tr> <td style="text-align: right;">+</td> <td>Authoritative sources in a hyperlinked environment</td> <td>(Paper) 49011</td> </tr> <tr> <td style="text-align: right;">+</td> <td>H.3.5</td> <td>(addresses-generic-area-of-interest) 15850</td> </tr> <tr> <td style="text-align: right;">+</td> <td>G.1.3</td> <td>(addresses-generic-area-of-interest) 8936</td> </tr> <tr> <td style="text-align: right;">+</td> <td>Local methods for estimating pagerank values</td> <td>(Paper) 6606</td> </tr> <tr> <td style="text-align: right;">+</td> <td>F.2.1</td> <td>(addresses-generic-area-of-interest) 1710</td> </tr> <tr> <td style="text-align: right;">+</td> <td>C.2.1</td> <td>(addresses-generic-area-of-interest) 516</td> </tr> <tr> <td style="text-align: right;">+</td> <td>H.3.1</td> <td>(addresses-gener</td> </tr> <tr> <td style="text-align: right;">+</td> <td>Ranking the web frontier</td> <td>(Paper)</td> </tr> <tr> <td style="text-align: right;">+</td> <td>Focused crawling</td> <td>(Paper)</td> </tr> </table>	+	H.3.3	(addresses-generic-area-of-interest) 2804692	+	The anatomy of a large-scale hypertextual Web search engine	(Paper) 122898	+	Authoritative sources in a hyperlinked environment	(Paper) 49011	+	H.3.5	(addresses-generic-area-of-interest) 15850	+	G.1.3	(addresses-generic-area-of-interest) 8936	+	Local methods for estimating pagerank values	(Paper) 6606	+	F.2.1	(addresses-generic-area-of-interest) 1710	+	C.2.1	(addresses-generic-area-of-interest) 516	+	H.3.1	(addresses-gener	+	Ranking the web frontier	(Paper)	+	Focused crawling	(Paper)
+	H.3.3	(addresses-generic-area-of-interest) 2804692																																
+	The anatomy of a large-scale hypertextual Web search engine	(Paper) 122898																																
+	Authoritative sources in a hyperlinked environment	(Paper) 49011																																
+	H.3.5	(addresses-generic-area-of-interest) 15850																																
+	G.1.3	(addresses-generic-area-of-interest) 8936																																
+	Local methods for estimating pagerank values	(Paper) 6606																																
+	F.2.1	(addresses-generic-area-of-interest) 1710																																
+	C.2.1	(addresses-generic-area-of-interest) 516																																
+	H.3.1	(addresses-gener																																
+	Ranking the web frontier	(Paper)																																
+	Focused crawling	(Paper)																																

[Graph info](#)
[Reindex](#)
[Reload Graph Data](#)
[Create Result Graph](#)
[Annotate](#)

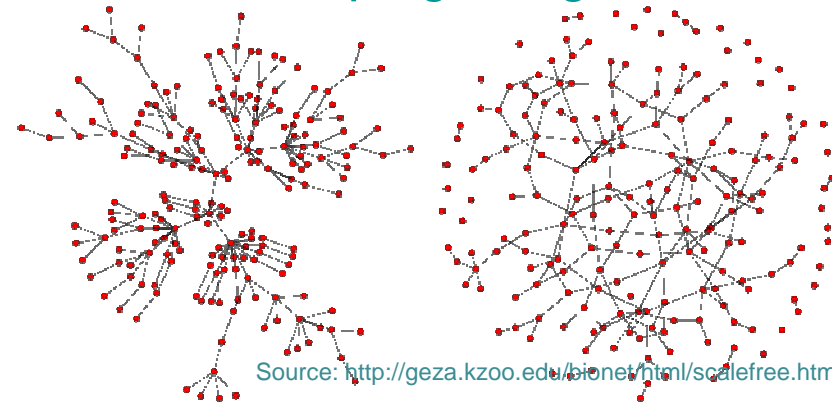
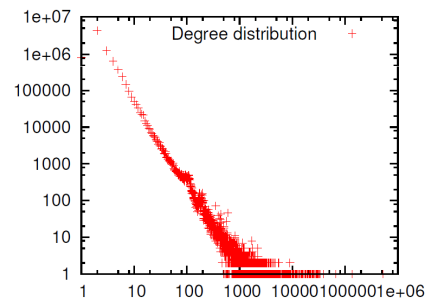


- Konverzia ACM LinkedData na jednoduchý graf pre gSemSearch
 - Experiment na hľadanie relácií a navigáciu
 - Pri konverzii na jednoduchší graf zanedbanie typov vzťahov: niekedy problém





- Storage for graphs
- Optimized for graph traversing and spread of activation
- Faster than Neo4j for graph traversing operations
- Supports Blueprints API
- <https://simplegdb.svn.sourceforge.net/svnroot/simplegdb/Sgdb3>



Source: <http://geza.kzoo.edu/bionet/html/scalefree.html>

- Graph Database Benchmarks

- Graph Traversal Benchmark for Graph Databases

- <http://ups.savba.sk/~marek/gbench.html>

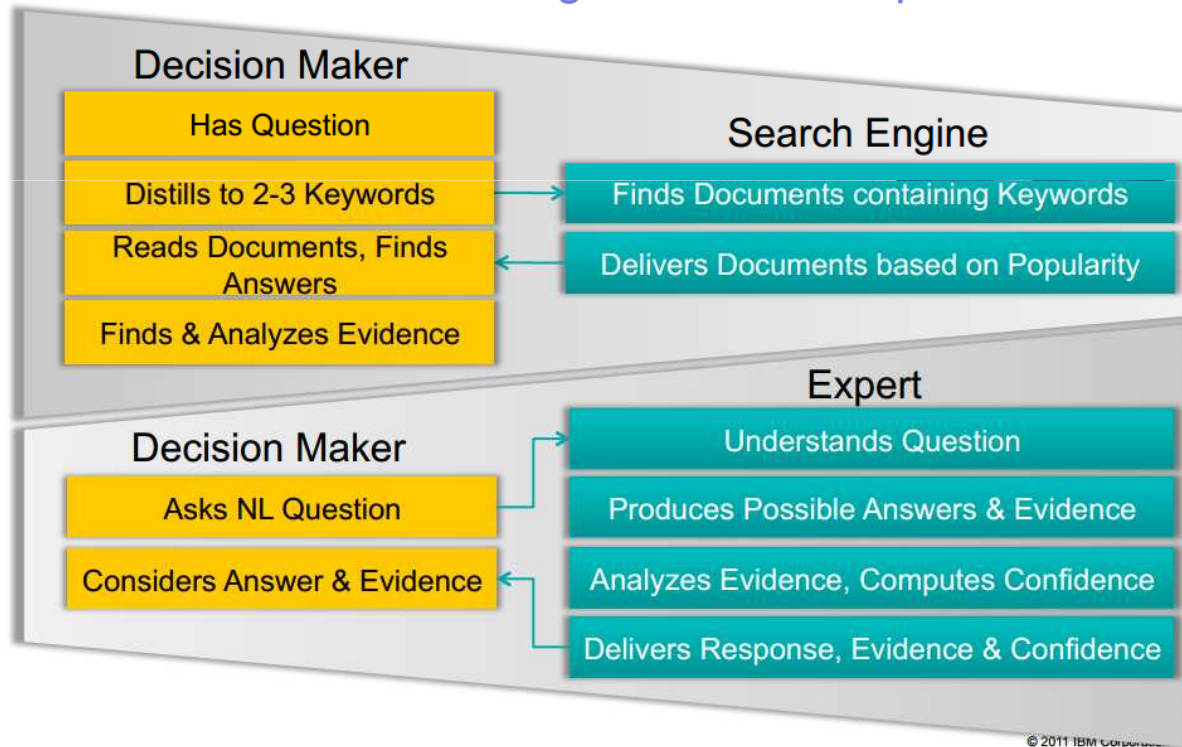
- Blueprints API - possibility to test compliant Graph databases





IBM Research

Informed Decision Making: Search vs. Expert Q&A





Európska únia
Európsky fond regionálneho rozvoja

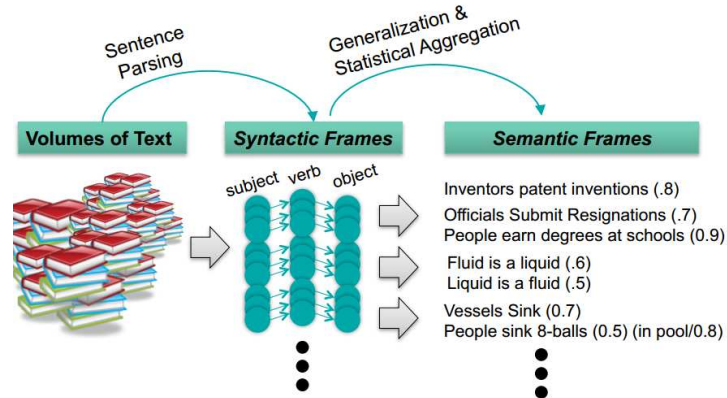
IBM Watson

[Perrone11]



IBM Research

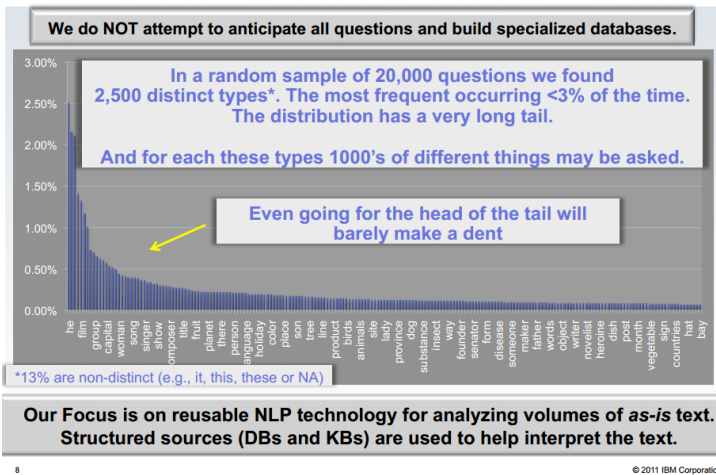
Automatic Learning From "Reading"



© 2011 IBM Corporation

IBM Research

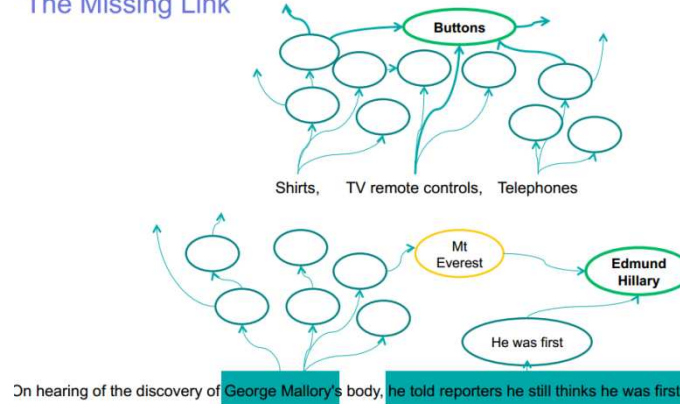
Broad Domain



© 2011 IBM Corporation

IBM Research

The Missing Link

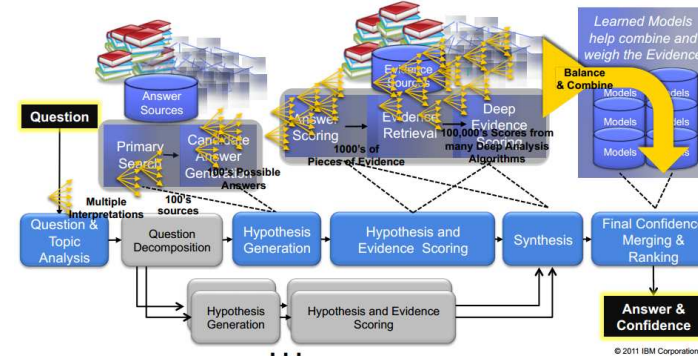


© 2011 IBM Corporation

IBM Research

DeepQA: The Technology Behind Watson

Massively Parallel Probabilistic Evidence-Based Architecture
Generates and scores many hypotheses using a combination of 1000's **Natural Language Processing, Information Retrieval, Machine Learning and Reasoning Algorithms.** These gather, evaluate, weigh and balance different types of **evidence** to deliver the answer with the best support it can find.



© 2011 IBM Corporation



- Wikipedia
 - 6 million articles
 - 40 GB text
- DBpedia
 - Triples
 - Good categories for articles
- Freebase
 - 170 GB triples
 - 40 million topics
 - 1.2 billion triples



- [Activity](#) (edit)
 - [Game](#) (edit)
 - [BoardGame](#) (edit)
 - [Sport](#) (edit)
 - [Athletics](#) (edit)
 - [Boxing](#) (edit)
 - [BoxingCategory](#) (edit)
 - [BoxingStyle](#) (edit)
 - [HorseRiding](#) (edit)
- [Agent](#) (edit)
 - [Deity](#) (edit)
 - [Family](#) (edit)
 - [Organisation](#) (edit)
 - [Band](#) (edit)
 - [Broadcaster](#) (edit)
 - [BroadcastNetwork](#) (edit)
 - [RadioStation](#) (edit)
 - [TelevisionStation](#) (edit)



```

ns:m.012rkqx ns:type.object.type ns:common.topic.
ns:m.012rkqx ns:type.object.name "High Fidelity"@en.
ns:m.012rkqx ns:type.object.type ns:music.single.
ns:m.012rkqx ns:type.object.key ns:authority.musicbrainz.name.TRACK3987054.
ns:m.012rkqx ns:type.object.type ns:music.recording.
ns:m.012rkqx key:authority.musicbrainz "258c45bd-4437-4580-8988-b3f3be975f9c".
ns:m.012rkqx key:authority.musicbrainz.name "TRACK3987054".
ns:m.012rkqx rdfs:label "High Fidelity"@en.
ns:m.012rkqx rdfs:type ns:common.topic.
ns:m.012rkqx rdfs:type ns:music.single.
ns:m.012rkqx rdfs:type ns:music.recording.
  
```


Ted Nugent

From Wikipedia, the free encyclopedia

"Nuge" redirects here. For the skateboarder, see Don Nguyen.

Theodore Anthony "Ted" Nugent (/tɛd nuːdʒɪnt/; born December 13, 1948) is an American rock musician from Detroit, Michigan. Nugent initially gained fame as the lead guitarist of The Amboy Dukes before embarking on a solo career. His hits, mostly coming in the 1970s, such as "Stranglehold", "Cat Scratch Fever", "Wango Tango", and "Great White Buffalo", as well as his 1960s Amboy Dukes hit "Journey to the Center of the Mind", remain popular today, and are played on classic rock and less frequently active rock radio stations. He is also noted for his staunch conservative political views and his strong advocacy of hunting and gun ownership rights, which some^[*who?*] have described as controversial.^{[1][2]}

Ted Nugent



Nugent in concert in Indianapolis, July 31, 2013.

Background information

Birth name Theodore Anthony Nugent

Also known as The Nuge, Motor City Madman, Uncle Ted

Born December 13, 1948 (age 64)
Redford, Michigan, U.S.

Genres Rock, hard rock

Occupations Musician

Instruments Guitar, vocals, bass guitar

Contents [hide]

- 1 Early life
- 2 Career
 - 2.1 Amboy Dukes
 - 2.2 Solo career
 - 2.3 Influences
 - 2.4 Damn Yankees



Smartphone /m/0169zh

Video Game Platform, Computing Platform, Literature Subject, Product category, Competitive Space

alias: smart phone, smartphone

A smartphone, or smart phone, is a mobile phone built on a mobile operating system, with more advanced computing capability and connectivity than a feature phone. The first smartphones combined the functions of a personal digital assistant with a... [Wikipedia]

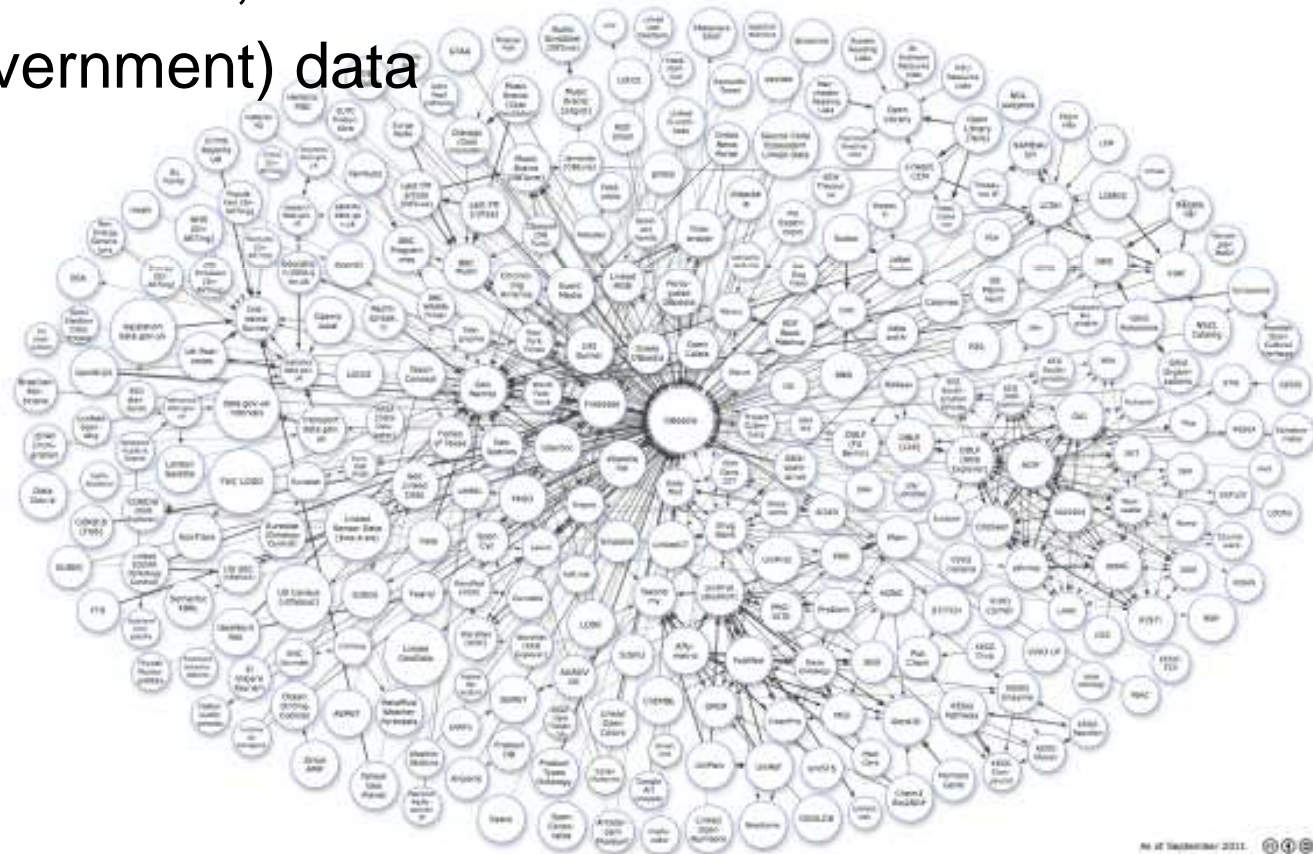


Európska únia
Európsky fond regionálneho rozvoja

Linked Data cloud



- Triples (graph)
- DBPedia, Geo, people, publications, medicine, ...
- EU public (government) data

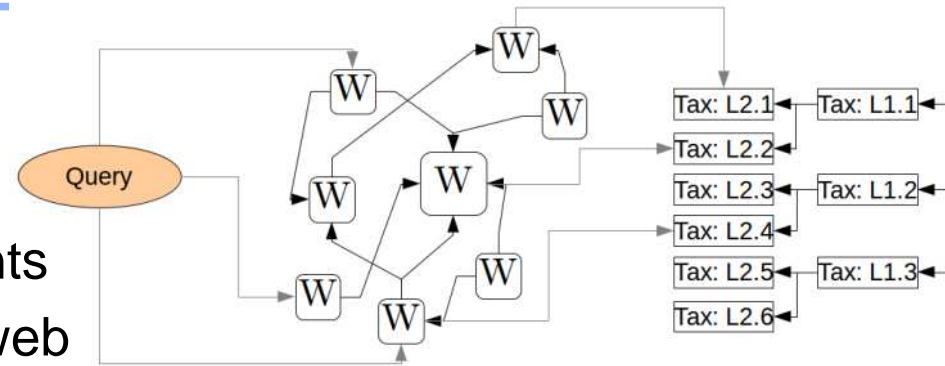


As of September 2011.

Query Categorization (QC)



- Usual approach for QC:
 - Get results for query
 - Categorize returned documents
- Best algorithms work with entire web (search API)



Query	Categories
apple	Computers \ Hardware Living \ Food & Cooking
FIFA 2006	Sports \ Soccer Sports \ Schedules & Tickets Entertainment \ Games & Toys
cheesecake recipes	Living \ Food & Cooking Information \ Arts & Humanities
friendships poem	Information \ Arts & Humanities Living \ Dating & Relationships

fifa 2006

Web Images Maps Shopping Applications More Search tools

About 95,700,000 results (0.30 seconds)

[2006 FIFA World Cup - Wikipedia, the free encyclopedia](#)
en.wikipedia.org/wiki/2006_FIFA_World_Cup
The 2006 FIFA World Cup was the 18th FIFA World Cup, the quadrennial international football world championship tournament. It was held from 9 June to 9 July ...
2006 FIFA World Cup Final - 2006 FIFA World Cup squads - Knockout stage - Petit

Images for **fifa 2006** - Report images



FIFA 06 - Download
fifa-06.en.softonic.com/
★★★★★ Rating: 4 - 520 votes - Free - Windows - Game
2006 FIFA World Cup - NBA Live 2001. EA SPORTS FIFA 06 brings the TOTAL FOOTBALL experience to your fingertips. It delivers a combination of attacking ...
Download - Clear filters - Related - See screenshots (7)

FIFA 2006 - Barcelona vs Real Madrid (El Clásico) - YouTube
www.youtube.com/watch?v=QIT221cFeZ4
Jan 30, 2012 - Uploaded by Lightning8S
GAME: FIFA 2006 DIFFICULTY: Professional HALF LENGTH: 6 Minutes STADIUM: Camp Nou PARTICIPANTS ...

2006 FIFA World Cup - Download
2006-fifa-world-cup-germany.en.softonic.com/
★★★★★ Rating: 4 - 257 votes - Free - Windows - Game
2006 FIFA World Cup, free download. 2006 FIFA World Cup Demo: Celebrate the passion of the World Cup. EA proves it just can't get enough with 2006 FIFA ...

2006 FIFA World Cup Germany™ - FIFA.com
www.fifa.com/worldcup/archive/germany2006/index.html
Jun 9, 2006 - A look back at the 2006 FIFA World Cup Germany™



- Each category represented by one or multiple wiki pages
- Example: *IAB19-6 Cell Phones* category
 - Mobile phone; Smartphone; Camera phone.

IAB1 Arts & Entertainment	The arts	Entertainment			
IAB1-1 Books & Literature	Book	Literature	Novel		
IAB1-2 Celebrity Fan/Gossip	Celebrity				
IAB1-3 Fine Art	Fine art				
IAB1-4 Humor	Humour	Fun			
IAB1-5 Movies	Film	Filmmaking	Film industry	Film genre	Movie theater
IAB1-6 Music	Music	Pop music	Classical music	Rock music	Music genre
IAB1-7 Television	Television program	Television	Serial (radio and television)		List of popular music genres
01-03 Performing art & Theatre	Performing arts	Performing arts	Theatre		
IAB2 Automotive	Automotive industry	Automobile	Car classification		
IAB2-1 Auto Parts	List of auto parts				
IAB2-2 Auto Repair	Automotive Service Excellence				
IAB2-3 Buying/Selling Cars	Cars	Used car	Car dealerships in North America		
IAB2-4 Car Culture	Effects of the automobile on societies				
IAB2-5 Certified Pre-Owned		Certified Pre-Owned			
IAB2-6 Convertible	Convertible				
IAB2-7 Coupe	Coupé	Compact car			
IAB2-8 Crossover	Crossover (automobile)				
IAB2-9 Diesel	Turbo-diesel				
IAB2-10 Electric Vehicle		Electric car			
IAB2-11 Hatchback	Hatchback				
IAB2-12 Hybrid	Hybrid vehicle				
IAB2-13 Luxury	Luxury vehicle				
IAB2-14 MiniVan	Minivan				
IAB2-15 Motorcycles	Motorcycle	List of motorcycle manufacturers			
IAB2-16 Off-Road Vehicles	Off-road vehicle				
IAB2-17 Performance Vehicles	Performance car				
IAB2-18 Pickup	Pickup truck				
IAB2-19 Road-Side Assistance	Roadside assistance	AAA (American Automobile Association)			
IAB2-20 Sedan	Sedan (automobile)				
IAB2-21 Trucks & Accessories	Truck	Truck accessory			
IAB2-22 Vintage Cars	Vintage car	History of the automobile			
IAB2-23 Wagon	Station wagon				
02-03 Campers & RVs	Truck camper	Recreational vehicle			
02-06 Concept Cars	Concept car				
02-08 SUV	Sport utility vehicle				

Mobile phone

From Wikipedia, the free encyclopedia

"Cell Phone" redirects here. For the film, see Cell Phone (film). For the Handphone film, see Handphone (film).

A **mobile phone** (also known as a **cellular phone**, **cell phone**, and a **hand phone**) is a device that can make and receive telephone calls over a radio link while moving around a wide geographic area. It does so by connecting to a cellular network provided by a mobile phone operator, allowing access to the public telephone network. By contrast, a cordless telephone is used only within the short range of a single, private base station.

In addition to telephony, modern mobile phones also support a wide variety of other services such as text messaging, MMS, email, Internet access, short-range wireless communications (infrared, Bluetooth), business applications, gaming and photography. Mobile phones that offer these and more general computing capabilities are referred to as smartphones.



The Qualcomm QCP-2700, a mid-1990s candybar style phone, and an iPhone 5, a current production smartphone.

Smartphone

From Wikipedia, the free encyclopedia

This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. (August 2012)

A **smartphone**, or **smart phone**, is a mobile phone built on a mobile operating system, with more advanced computing capability and connectivity than a feature phone (FPH). The first smartphones combined the functions of a personal digital assistant (PDA) with a mobile phone. Later models added the functionality of portable media players, low-end compact digital cameras, pocket video cameras, and GPS navigation units to form one multi-use device. Many modern smartphones also include high-resolution touchscreens and web browsers that display standard web pages as well as mobile-optimized sites. High-speed data access is provided by Wi-Fi and mobile broadband. In recent years, the rapid development of mobile app markets and of mobile commerce have been drivers of smartphone adoption.

The mobile operating systems (OS) used by modern smartphones include Google's Android, Apple's iOS, Nokia's Symbian, BlackBerry Ltd's BlackBerry OS, Samsung's Bada, Microsoft's Windows Phone, Hewlett-Packard's webOS, and embedded Linux distributions such as MeeGo and Meego. Such operating systems can be installed on many different phone models, and typically each device can receive multiple OS software updates over its lifetime. A few other upcoming operating systems are Mozilla's Firefox OS, Canonical Ltd's Ubuntu Phone, and Tizen.

Camera phone

From Wikipedia, the free encyclopedia

For the song by The Game, see Camera Phone (song).

See also: *Mobile phone* and *Videophone*

This article's factual accuracy may be compromised due to out-of-date information. Please update this article to reflect recent events or newly available information. (May 2012)



A **camera phone** is a mobile phone which is able to capture still photographs (and usually video). Since early in the 21st century the majority of mobile phones in use are camera phones.^[1]



- Seed concepts – manually assigned for each category
- Exploit Wiki link graph
- For all neighboring Wiki Concepts (to seed concepts)
 - Filter out concepts of type: Person, Location, Org, Work
 - Compute cosine similarity of adjacency lists
 - Extend concept names by alternative names (redirect pages from Wiki)

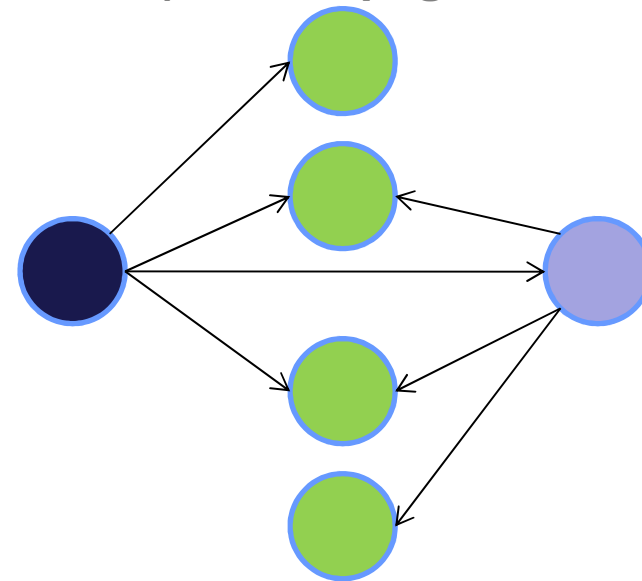
Smartphone

From Wikipedia, the free encyclopedia



This article **needs additional citations for verification**. Please help [improve this article](#) by adding citations to reliable sources. Unsourced material may be [challenged](#) and removed. (August 2013)

A **smartphone**, or **smart phone**, is a mobile phone built on a [mobile operating system](#) with more advanced computing capability and connectivity than a feature phone. ^{[[cite\]](#)]} The first smartphones combined the functions of a [personal digital assistant \(PDA\)](#) [Mobile operating system](#)





Európska únia
Európsky fond regionálneho rozvoja

Literatúra



- [Ulanoff] Lance Ulanoff: **Google Knowledge Graph Could Change Search Forever** <http://mashable.com/2012/02/13/google-knowledge-graph-change-search/>, 2012
- [facebook13] Sean Gallagher, **Knowing the score: How Facebook's Graph Search knows what you want**, <http://arstechnica.com/information-technology/2013/03/knowning-the-score-how-facebooks-graph-search-knows-what-you-want/>, 2013
- [Perrone11] Michael Perrone: **What is Watson – An Overview**, 2011, <http://static.usenix.org/event/lisa11/tech/slides/perrone.pdf>
- [WatsonJr] Tony Pearson: **IBM Watson - How to build your own "Watson Jr." in your basement**, 2012, https://www.ibm.com/developerworks/mydeveloperworks/blogs/InsideSystemStorage/entry/ibm_watson_how_to_build_your_own_watson_jr_in_your_basement7?lang=en
- [OpenNLP] OpenNLP: <http://www.slideshare.net/gagan1667/opennlp-demo>
- [TamingText] Ingersoll, G., Morton, T., & Farris, L. (2012). Taming Text: How to find, organize and manipulate it.
- [Zaragoza] Hugo Zaragoza. **Machine Learning and Information Retrieval**, ESSIR 2009 Lecture
- [Nigam] Kamal Nigam: **Generative Models for Text Classification and Information Extraction**, <http://www.cs.cmu.edu/~knigam/15-505/ie-lecture.ppt>



- [SemSets] CIGLAN, Marek - NoRVaG, Kjetil - HLUCHÝ, Ladislav. **The SenSets model for ad-hoc semantic list search.** In WWW'12 Proceedings of the 21st International Conference on World Wide Web. - New York : ACM, 2012, p. 131-140. ISBN 978-1-4503-1229-5. SCOPUS, <http://www2012.wwwconference.org/proceedings/proceedings/p131.pdf>
- [gSemSearch] LACLAVÍK, Michal - DLUGOLINSKÝ, Štefan - ŠELENG, Martin - CIGLAN, Marek - HLUCHÝ, Ladislav. **Emails as graph: relation discovery in email archive.** In WWW'12 Companion Proceedings of the 21st International Conference companion on World Wide Web. - New York : ACM, 2012, 841-846. ISBN 978-1-4503-1230-1. <http://www2012.wwwconference.org/proceedings/companion/p841.pdf> . SCOPUS
- [gBench] CIGLAN, Marek - AVERBUCH, Alex - HLUCHÝ, Ladislav. **Benchmarking traversal operations over graph databases.** In 2012 IEEE 28th International Conference on Data Engineering Workshops : proceedings. - Los Alamitos : IEEE Computer Society, 2012, p. 186-189. ISBN 978-1-4673-1640-8. SCOPUS
- [ontea_email] LACLAVÍK, Michal - DLUGOLINSKÝ, Štefan - ŠELENG, Martin - KVASSAY, Marcel - GATIAL, Emil - BALOGH, Zoltán - HLUCHÝ, Ladislav. **Email analysis and information extraction for enterprise benefit.** In **Computing and informatics**, 2011, vol. 30, no. 1, p. 57-87. (0.356 - IF2010). ISSN 0232-0274.
- [uiWeb] Dlugolinský, Štefan - Šeleng, Martin - Laclavík, Michal - Hluchý, Ladislav. **Distributed Web-scale Infrastructure for Crawling, Indexing and Search with Semantic Support.** In **Computer Science Journal**, 13 (4)



Európska únia
Európsky fond regionálneho rozvoja

Ďakujeme za pozornosť

